



UNIVERSITÉ
DE GENÈVE

FACULTY OF SCIENCE
Department of Informatics

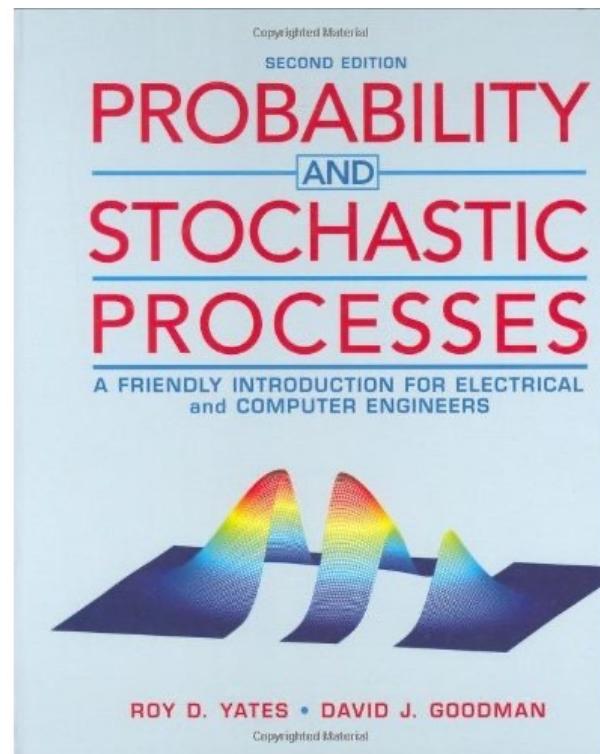
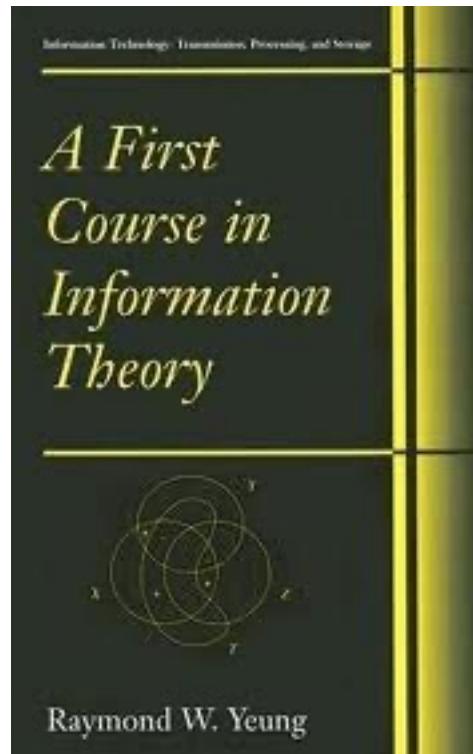
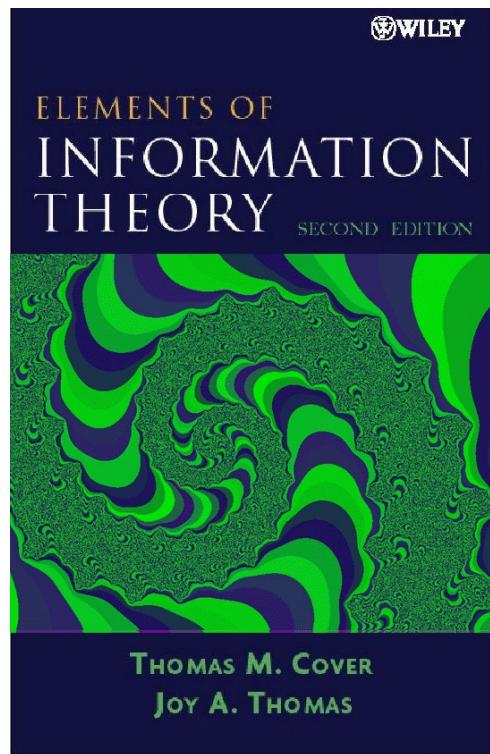


Information Theory for Data Science and Machine Learning

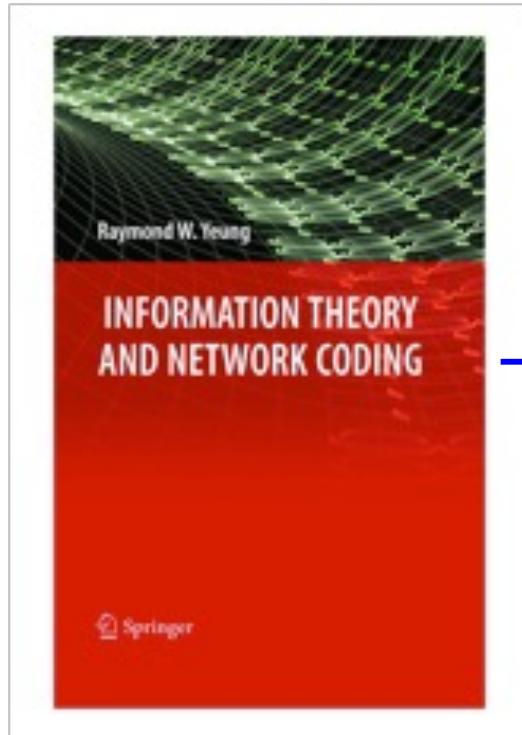
Slava Voloshynovskiy

<http://sip.unige.ch/education/information-theory/>

Recommended books



Recommended books



Free download as

Raymond W. Yeung

Information Theory and
Network Coding

SPIN Springer's internal project number, if known

January 28, 2008

<http://iest2.ie.cuhk.edu.hk/~whyeung/book2/>

<http://iest2.ie.cuhk.edu.hk/~whyeung/post/manuscript/main2.pdf>

Copyright notice

This course uses in part some materials from previously mentioned recommended sources.

The usage of the slides in commercial or educational purposes
is prohibited without authorization of the course leaders and
permission of the above document authors and copyright owners.

Content of course

- Theme 1 – Basic statistical data models
- Theme 2 – Information theoretic measures

Scope

- Information theoretic measures for discrete/continuous variables
 - Entropy
 - Conditional entropy
 - Joint entropy
 - Relative entropy (KL-divergence)
 - Cross entropy
 - Mutual information
 - Additional topics:
 - f-divergence
 - variational inference
 - practical computations from samples

Part I: discrete random variables and vectors

Part II: continuous random variables and vectors

Scope



Part I: discrete random variables and vectors

Scope

- Information theoretic measures for discrete/continuous variables

- Entropy
- Conditional entropy
- Joint entropy
- Relative entropy (KL-divergence)
- Cross entropy
- Mutual information
- Additional topics:

- f-divergence
- variational inference
- practical computations from samples

Entropy

Definition (entropy): *Entropy* of discrete r.v. $X \in \mathcal{X}$ with $\mathcal{X} = \{x_1, x_2, \dots, x_n\}$

$$H(X) = -\sum_{x \in \mathcal{X}} p_X(x) \log_2 p_X(x) = -E_{p_X}[\log p_X(x)] \quad (\text{bits})$$

$$H(X) = -\sum_{i=1}^n \Pr\{X = x_i\} \log_2 \Pr\{X = x_i\} \quad E_{p_X}[g(X)] = \sum_{x \in \mathcal{X}} p_X(x) g(x)$$

- When the base of logarithm is α , we use the notation $H_\alpha(X)$
- Entropy measures the **uncertainty** or **randomness** of a discrete random variable
- The units of entropy
 - bit – if $\alpha = 2$
 - nat – if $\alpha = e$
- $H(X)$ depends only on the distribution of X but **not** on the actual value taken by X , hence also it is reflected as $H(p)$ or $H(p(x))$

Entropy

- Examples

- Entropy of binary random variable with the equiprobable states

$$H(X) = -\left[\frac{1}{2} \log_2 \frac{1}{2} + \frac{1}{2} \log_2 \frac{1}{2} \right] = \log_2 2 = 1 \text{ bit}$$

- Entropy of a random variable with 27 (a model of alphabet) equiprobable states

$$H(X) = -\frac{1}{27} \sum_{i=1}^{27} \log_2 \frac{1}{27} = \log_2 27 = 4.7549 \text{ bit}$$

- Note: in reality $H(X) = 3.98 \text{ bit}$
 - Entropy of a random variable with N equiprobable states/symbols

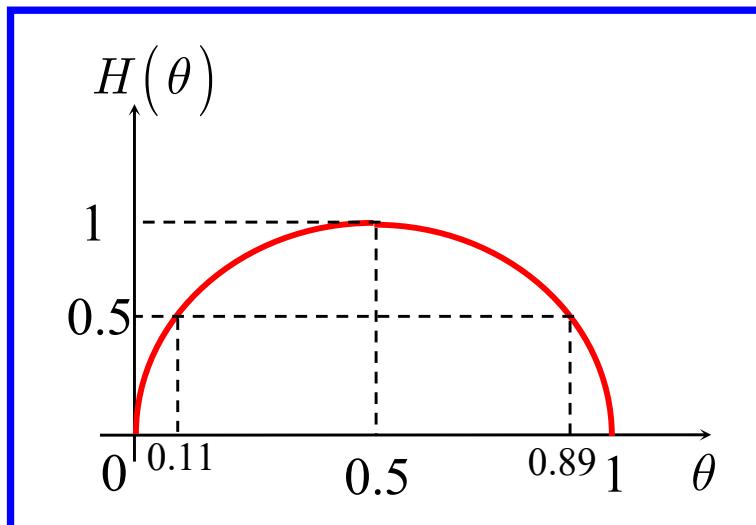
$$H(X) = -\frac{1}{N} \sum_{i=1}^N \log_2 \frac{1}{N} = \log_2 N \text{ bit}$$

Entropie: binary random variable

- Binary random variable $X \in \mathcal{X} = \{0,1\}$

$$\Pr\{X = 0\} = 1 - \theta \quad \Rightarrow \Pr[X = x] = \theta^x (1 - \theta)^{1-x}$$
$$\Pr\{X = 1\} = \theta$$

$$H(X) = -\theta \log_2 \theta - (1 - \theta) \log_2 (1 - \theta) \coloneqq H(\theta) \text{ or } H_2(\theta)$$



The function $H(\theta)$ is:

- Symmetric wrt $p = 0.5$
- Maximum ($H(p) = 1$) at $\theta = 0.5$

The entropy is maximal, if the symbols are equiprobable.

Joint entropy

Definition (joint entropy): *Joint entropy* of two discrete random variables X and Y is defined by:

$$H(X, Y) = - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{X,Y}(x, y) \log_2 p_{X,Y}(x, y) = -E_{p_{X,Y}} [\log p_{X,Y}(x, y)]$$

$$\mathcal{X} = \{x_1, x_2, \dots, x_n\} \quad \mathcal{Y} = \{y_1, y_2, \dots, y_m\}$$

$$H(X, Y) = H(Y, X)$$

Conditional entropy

Definition (conditional entropy): *Conditional entropy* of r.v. X given Y is defined by:

$$H(X|Y) = \sum_{y \in \mathcal{Y}} p_Y(y) H(X|Y=y) = -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{X,Y}(x,y) \log p_{X|Y}(x|y)$$

$$\mathcal{X} = \{x_1, x_2, \dots, x_n\} \quad \mathcal{Y} = \{y_1, y_2, \dots, y_m\}$$

$$H(X|Y) \neq H(Y|X)$$

$$H(X|Y=y) = -\sum_{i=1}^n p_{X|Y}(x_i|y) \log p_{X|Y}(x_i|y)$$

$$H(X|Y) = -\sum_{j=1}^m p_Y(y_j) \left[\sum_{i=1}^n p_{X|Y}(x_i|y_j) \log p_{X|Y}(x_i|y_j) \right] =$$

$$= -\sum_{i=1}^n \sum_{j=1}^m p_{X,Y}(x_i, y_j) \log p_{X|Y}(x_i|y_j)$$

Entropy: computation

Exemple:

- (a) $H(X) - ?$
- (b) $H(Y) - ?$
- (c) $H(X|Y) - ?$
- (d) $H(Y|X) - ?$
- (e) $H(X, Y) - ?$

$\backslash X$	1	2	3	4
1	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{32}$	$\frac{1}{32}$
2	$\frac{1}{16}$	$\frac{1}{8}$	$\frac{1}{32}$	$\frac{1}{32}$
3	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{16}$
4	$\frac{1}{4}$	0	0	0

Entropy: computation

Example: basic entropy computation

	X	Y	
1	$p(X = 1, Y = 1) p(X = 2, Y = 1) p(X = 3, Y = 1) p(X = 4, Y = 1)$		$p(Y = 1)$
2	$p(X = 1, Y = 2) p(X = 2, Y = 2) p(X = 3, Y = 2) p(X = 4, Y = 2)$		$p(Y = 2)$
3	$p(X = 1, Y = 3) p(X = 2, Y = 3) p(X = 3, Y = 3) p(X = 4, Y = 3)$		$p(Y = 3)$
4	$p(X = 1, Y = 4) p(X = 2, Y = 4) p(X = 3, Y = 4) p(X = 4, Y = 4)$		$p(Y = 4)$
	$p(X = 1)$	$p(X = 2)$	$p(X = 3)$
			$p(X = 4)$

Entropy: computation

Example: basic entropy computation

$X \setminus Y$	1	2	3	4	$p(Y = i)$
$p(X = j)$	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{8}$	$\frac{1}{8}$	
1	$\frac{1}{8}$	$\frac{1}{16}$	$\frac{1}{32}$	$\frac{1}{32}$	$p(Y = 1) = \frac{1}{4} \quad p(Y = 1) = \sum_{j=1}^4 p(X = j, Y = 1)$
2	$\frac{1}{16}$	$\frac{1}{8}$	$\frac{1}{32}$	$\frac{1}{32}$	$p(Y = 2) = \frac{1}{4} \quad p(Y = 2) = \sum_{j=1}^4 p(X = j, Y = 2)$
3	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{16}$	$\frac{1}{16}$	$p(Y = 3) = \frac{1}{4} \quad p(Y = 3) = \sum_{j=1}^4 p(X = j, Y = 3)$
4	$\frac{1}{4}$	0	0	0	$p(Y = 4) = \frac{1}{4} \quad p(Y = 4) = \sum_{j=1}^4 p(X = j, Y = 4)$

Entropy: computation

$$(a) \quad H(X) = -\sum_{j=1}^4 p(X=j) \log p(X=j)$$

$$p(X=j) = \left\{ \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8} \right\}$$

$$\begin{aligned} H(X) &= -\left[\frac{1}{2} \log 2^{-1} + \frac{1}{4} \log 2^{-2} + \frac{1}{8} \log 2^{-3} + \frac{1}{8} \log 2^{-3} \right] = \\ &= \frac{1}{2} + \frac{2}{4} + \frac{3}{8} + \frac{3}{8} = \frac{7}{4} \text{ bit.} \end{aligned}$$

Entropy: computation

$$(b) \quad H(Y) = -\sum_{i=1}^4 p(Y=i) \log p(Y=i)$$

$$p(Y = i) = \left\{ \frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4} \right\}$$

$$\begin{aligned} H(Y) &= -\left[\frac{1}{4} \log 2^{-2} + \frac{1}{4} \log 2^{-2} + \frac{1}{4} \log 2^{-2} + \frac{1}{4} \log 2^{-2} \right] = \\ &= \frac{1}{2} + \frac{1}{2} + \frac{1}{2} + \frac{1}{2} = 2 \text{ bit.} \end{aligned}$$

Entropy: computation

$$\begin{aligned}
 (c) \quad H(X|Y) &= \sum_{i=1}^4 p(Y=i)H(X|Y=i) = \\
 &= p(Y=1)H(X|Y=1) + p(Y=2)H(X|Y=2) + p(Y=3)H(X|Y=3) \\
 &\quad + p(Y=4)H(X|Y=4)
 \end{aligned}$$

$$p(Y=i) = \left\{ \frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4} \right\}$$

$$H(X|Y) = \frac{1}{4}H(X|Y=1) + \frac{1}{4}H(X|Y=2) + \frac{1}{4}H(X|Y=3) + \frac{1}{4}H(X|Y=4)$$

Entropy: computation

$$\underline{H(X|Y=1)} = -\sum_{i=1}^4 p(X=i|Y=1) \log p(X=i|Y=1)$$

Recall:

$$p(x,y) = p(x|y)p(y) \Rightarrow p(x|y) = \frac{p(x,y)}{p(y)}$$

$$p(X=1|Y=1) = \frac{p(X=1, Y=1)}{p(Y=1)} = \frac{1/8}{1/4} = \frac{1}{2}$$

From table

$$p(X=2|Y=1) = \frac{p(X=2, Y=1)}{p(Y=1)} = \frac{1/16}{1/4} = \frac{1}{4}$$

$$p(Y=i) = \left\{ \frac{1}{4}, \frac{1}{4}, \frac{1}{4}, \frac{1}{4} \right\}$$

$$p(X=3|Y=1) = \frac{p(X=3, Y=1)}{p(Y=1)} = \frac{1/32}{1/4} = \frac{1}{8}$$

$$p(X=4|Y=1) = \frac{p(X=4, Y=1)}{p(Y=1)} = \frac{1/32}{1/4} = \frac{1}{8}$$

Entropy: computation

$$\begin{aligned}
 H(X|Y=1) &= -\sum_{i=1}^4 p(X=i|Y=1) \log p(X=i|Y=1) = \\
 &= -\left[\frac{1}{2} \log 2^{-1} + \frac{1}{4} \log 2^{-2} + \frac{1}{8} \log 2^{-3} + \frac{1}{8} \log 2^{-3} \right] = \\
 &= \frac{1}{2} + \frac{1}{2} + \frac{3}{8} + \frac{3}{8} = \frac{7}{4}.
 \end{aligned}$$

$$H(X|Y=2) = \frac{7}{4}, \quad H(X|Y=3) = 2, \quad H(X|Y=2) = 0.$$

$$\begin{aligned}
 H(X|Y) &= \frac{1}{4} H(X|Y=1) + \frac{1}{4} H(X|Y=2) + \frac{1}{4} H(X|Y=3) + \frac{1}{4} H(X|Y=4) \\
 &= \frac{1}{4} \cdot \frac{7}{4} + \frac{1}{4} \cdot \frac{7}{4} + \frac{1}{4} \cdot 2 + \frac{1}{4} \cdot 0 = \frac{11}{8} \text{ bit}
 \end{aligned}$$

Entropy: computation

$$\begin{aligned}
 (d) \quad H(Y|X) &= \sum_{j=1}^4 p(X=j)H(Y|X=j) = \\
 &= p(X=1)H(Y|X=1) + p(X=2)H(Y|X=2) + p(X=3)H(Y|X=3) \\
 &\quad + p(X=4)H(Y|X=4)
 \end{aligned}$$

$$p(X = j) = \left\{ \frac{1}{2}, \frac{1}{4}, \frac{1}{8}, \frac{1}{8} \right\}$$

$$H(Y|X) = \frac{1}{2}H(Y|X=1) + \frac{1}{4}H(Y|X=2) + \frac{1}{8}H(Y|X=3) + \frac{1}{8}H(Y|X=4)$$

You can finalize it yourselves...

Entropy: computation

$$\Rightarrow H(Y|X) = \frac{13}{8} \text{ bit.}$$

$$H(X|Y) \neq H(Y|X)$$

Entropy: computation

$$\begin{aligned}
 (e) \quad H(X, Y) &= \sum_{i=1}^4 \sum_{j=1}^4 p(X = j, Y = i) \log p(X = j, Y = i) \\
 &= - \left[2 \times \frac{1}{8} \log 2^{-3} + 6 \times \frac{1}{16} \log 2^{-4} + 4 \times \frac{1}{32} \log 2^{-5} + \frac{1}{4} \log 2^{-2} \right] = \\
 &= \frac{3}{4} + \frac{6}{4} + \frac{5}{8} + \frac{2}{4} = \frac{27}{8} \text{ bit.}
 \end{aligned}$$

$$H(X, Y) = H(X) + H(Y | X)$$



$\frac{27}{8} \text{ bit.}$
 $\frac{7}{4} + \frac{13}{8} = \frac{27}{8} \text{ bit.}$

Summary: properties of entropy

- Property 1: non-negativity of entropy
- Property 2: chain rule for entropy
- Property 3: increase of entropy by “adding” new randomness
- Property 4: decrease of entropy by conditioning
- Property 5: uniform pmf has the largest entropy

Summary: properties of entropy

- Property 1: entropy is strictly positive

$$H(X) \geq 0$$

$$H(X, Y) \geq 0$$

$$H(X|Y) \geq 0$$

$$0 \leq p_X(x) \leq 1$$

$$0 \leq p_{X,Y}(x, y) \leq 1 \Rightarrow \log p_{X,Y}(x, y) \leq 0$$

$$0 \leq p_{X|Y}(x|y) \leq 1 \quad \log p_{X|Y}(x|y) \leq 0$$

$$0 \leq -\sum p_X(x) \log p_X(x)$$

$$\Rightarrow 0 \leq -\sum p_{X,Y}(x, y) \log p_{X,Y}(x, y)$$

$$0 \leq -\sum p_{X|Y}(x|y) \log p_{X|Y}(x|y)$$

Summary: properties of entropy

- Property 2: chain rule for entropy (case of 2 variables)

$$p(x,y) = p(x|y)p(y) = p(y|x)p(x) \Leftrightarrow p_{X,Y}(x,y) = p_{X|Y}(x|y)p_Y(y) = p_{Y|X}(y|x)p_X(x)$$

$$H(X,Y) = H(X|Y) + H(Y) = H(Y|X) + H(X)$$

Proof

$$-\log p_{X,Y}(x,y) = -\log p_{Y|X}(y|x)p_X(x) = -\log p_{Y|X}(y|x) - \log p_X(x).$$

$$\begin{aligned} H(X,Y) &= E_{p_{X,Y}}[-\log p_{X,Y}(x,y)] = E_{p_{X,Y}}[-\log p_X(x)] + E_{p_{X,Y}}[-\log p_{Y|X}(y|x)] \\ &= H(X) + H(Y|X). \end{aligned}$$

Summary: properties of entropy

- Property 2: chain rule for entropy (case of N variables)

$$p(x_1, \dots, x_N) = p(x_1) \prod_{i=2}^N p(x_i | x_{i-1}, \dots, x_1)$$

$$H(X_1, X_2, \dots, X_N) = \sum_{i=1}^N H(X_i | X_{i-1}, \dots, X_1)$$

Proof

$$H(X_1, X_2) = H(X_1) + H(X_2 | X_1) \quad \Rightarrow \quad H(X_1, X_2) \stackrel{-=, \text{ if independent}}{\leq} H(X_1) + H(X_2)$$

$$H(X_1, X_2, X_3) = H(X_1) + H(X_2, X_3 | X_1) = H(X_1) + H(X_2 | X_1) + H(X_3 | X_2, X_1)$$

$$\begin{aligned} H(X_1, X_2, \dots, X_N) &= H(X_1) + H(X_2 | X_1) + \dots + H(X_N | X_{N-1}, \dots, X_1) = \\ &= \sum_{i=1}^N H(X_i | X_{i-1}, \dots, X_1) \end{aligned}$$

Summary: properties of entropy

- Property 2: chain rule for entropy (case of N variables)

$$H(X_1, X_2, \dots, X_N) = \sum_{i=1}^N H(X_i \mid \underbrace{X_{i-1}, \dots, X_1}_{})$$



$$\mathbf{X} \triangleq \begin{bmatrix} | \\ \mathbf{x} \\ | \end{bmatrix} = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_N \end{bmatrix}$$

$$p_{\mathbf{x}}(\mathbf{x}) = p_{\mathbf{x}}(x_1, x_2, \dots, x_N)$$

$$\mathbf{X} \triangleq \begin{bmatrix} | \\ \mathbf{x} \\ | \end{bmatrix} = \begin{bmatrix} \vdots \\ X_{i-2} \\ X_{i-1} \\ X_i \\ \vdots \end{bmatrix}$$
$$\prod_{i=2}^N p(x_i \mid \mathbf{x}_{i-1}, \mathbf{x}_{i-2})$$

Summary: properties of entropy

- Property 2: chain rule for entropy (case of N variables)

$$H(X_1, X_2, \dots, X_N) \leq \sum_{i=1}^N H(X_i)$$

equality for independent $\{X_i\}$

Proof

For Independent $\{X_i\}$

$$H(X_2 | X_1) = H(X_2)$$

$$H(X_3 | X_2, X_1) = H(X_3)$$

$$H(X_i | X_{i-1}, \dots, X_1) = H(X_i)$$

$$H(X_1, X_2, \dots, X_N) = H(X_1) + H(X_2 | X_1) + \dots + H(X_N | X_{N-1}, \dots, X_1) =$$

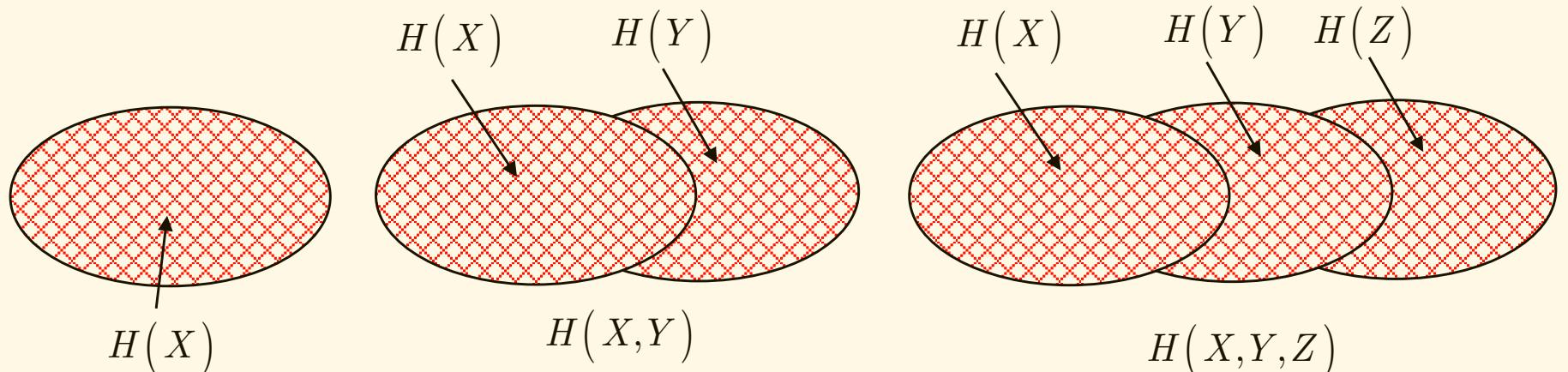
$$H(X_1, X_2, \dots, X_N) = H(X_1) + H(X_2) + \dots + H(X_N) = \sum_{i=1}^N H(X_i)$$

Summary: properties of entropy

- Property 3: increase of entropy by adding new randomness

$$H(X) \leq H(X, Y) \leq H(X, Y, Z) \leq \dots$$

$$H(X) = H(X, Y) \Leftrightarrow Y = f(X) \quad H(X, Y) = H(X) + \underbrace{H(Y|X)}_0$$
$$H(X, Y) = H(X, Y, Z) \Leftrightarrow Z = f(X, Y)$$



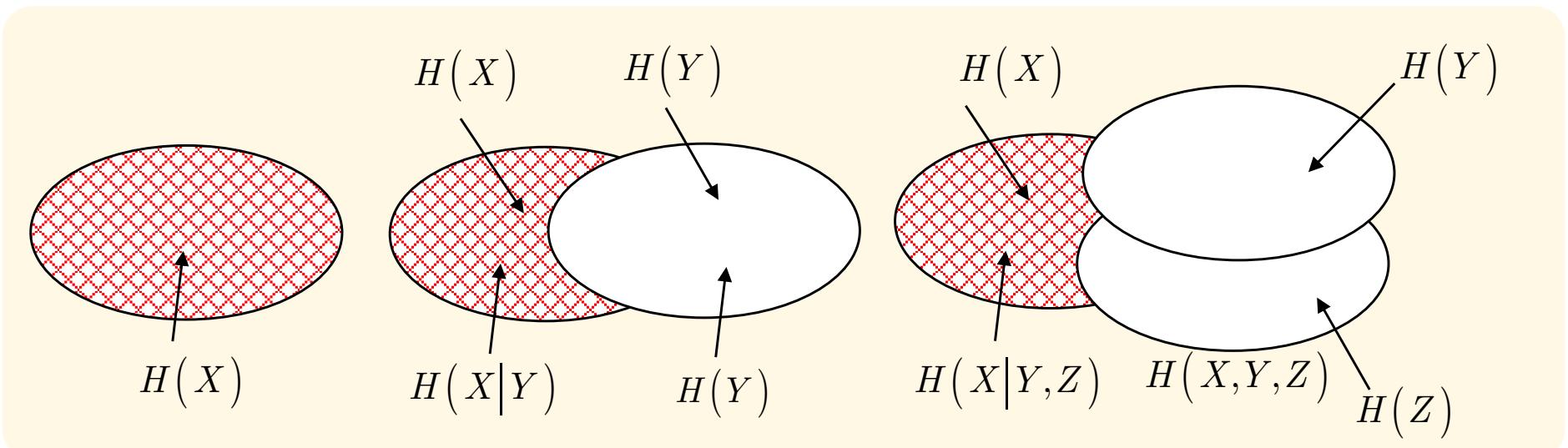
Summary: properties of entropy

- Property 4: conditioning reduces entropy

$$H(X) \geq H(X|Y) \geq H(X|Y,Z) \geq \dots$$

$$H(X) = H(X|Y) \Leftrightarrow X \perp Y$$

$$H(X|Y) = H(X|Y,Z) \Leftrightarrow X \perp Z|Y \dots$$



Summary: properties of entropy

- Property 5: uniform pmf has the largest entropy

$$H(X) \leq \log_2 |\mathcal{X}|$$

|
equality, iff $\Pr[X = x_i] = \frac{1}{N} = \frac{1}{|\mathcal{X}|}$

- Proof

Let $\mathcal{X} = \{x_1, x_2, \dots, x_N\}$ and $\Pr[X = x_i] = p_i$

$$\begin{aligned} H(X) &= -\sum_{x \in \mathcal{X}} \Pr[X = x] \log_2 \Pr[X = x] = \\ &= -(p_1 \log_2 p_1 + p_2 \log_2 p_2 + \dots + p_N \log_2 p_N) \end{aligned}$$

- Our goal: to find a distribution $\{p_i\}_{i=1}^N$ that has the maximum entropy $H(X)$

Summary: properties of entropy

- Recall: constrained optimization

$$\left(\hat{x}_1, \hat{x}_2, \dots, \hat{x}_N \right) = \arg \min_{x_1, x_2, \dots, x_N: \text{ such that } g(x_1, x_2, \dots, x_N) = k} f(x_1, x_2, \dots, x_N)$$

- We will use the method of Lagrange multipliers (unconstrained regularization)

$$\mathcal{L}(x_1, x_2, \dots, x_N) = f(x_1, x_2, \dots, x_N) + \lambda(g(x_1, x_2, \dots, x_N) - k)$$

$$\left(\hat{x}_1, \hat{x}_2, \dots, \hat{x}_N \right) = \arg \min_{x_1, x_2, \dots, x_N} \mathcal{L}(x_1, x_2, \dots, x_N)$$

- Solution $\frac{\partial \mathcal{L}(x_1, x_2, \dots, x_N)}{\partial x_i} = \frac{\partial f(x_1, \dots, x_i, \dots, x_N)}{\partial x_i} + \lambda \frac{\partial g(x_1, \dots, x_i, \dots, x_N)}{\partial x_i} = 0$

$$\frac{d\mathcal{L}(x_1, x_2, \dots, x_N)}{d\lambda} = g(x_1, x_2, \dots, x_N) - k = 0$$

Summary: properties of entropy

- In our case

$$f(p_1, p_2, \dots, p_N) = -\sum_{i=1}^N p_i \log_2 p_i \quad g(x_1, x_2, \dots, x_N) = \sum_{i=1}^N p_i = 1$$

- The Lagrange cost

$$\mathcal{L}(x_1, x_2, \dots, x_N) = f(x_1, x_2, \dots, x_N) + \lambda(g(x_1, x_2, \dots, x_N) - k)$$

$$\mathcal{L}(p_1, p_2, \dots, p_N) = -(p_1 \log_2 p_1 + p_2 \log_2 p_2 + \dots + p_N \log_2 p_N) + \lambda(p_1 + p_2 + \dots + p_N - 1)$$

- Solution

$$\begin{cases} \frac{\partial \mathcal{L}(p_1, p_2, \dots, p_N)}{\partial p_i} = -\log_2 p_i - 1 + \lambda = 0; \\ \frac{\partial \mathcal{L}(p_1, p_2, \dots, p_N)}{\partial p_j} = -\log_2 p_j - 1 + \lambda = 0; \quad \Rightarrow \quad -\log_2 p_i - 1 + \lambda = -\log_2 p_j - 1 + \lambda; \\ \frac{\partial \mathcal{L}(p_1, p_2, \dots, p_N)}{\partial \lambda} = p_1 + p_2 + \dots + p_N - 1 = 0; \end{cases}$$
$$\log_2 p_i = \log_2 p_j \quad \Rightarrow \quad p_i = p_j = p = \frac{1}{N}$$

$$p_i = \frac{1}{|\mathcal{X}|}; H_{\max}(X) = -\sum_{x \in \mathcal{X}} p(x) \log_2 p(x) = \log_2 |\mathcal{X}|$$

Relative entropy

Definition (relative entropy): *Relative entropy* or *Kullback-Leibler divergence (KLD)* between pmfs $p(x)$ and $q(x)$:

$$D_{\text{KL}}(p \parallel q) = D(p \parallel q) = \sum_{x \in \mathcal{X}} p(x) \log_2 \frac{p(x)}{q(x)} = E_{p(x)} \left[\log_2 \frac{p(x)}{q(x)} \right]$$

with the conventions: $0 \log_2 \frac{0}{q(x)} = 0, \forall q(x)$

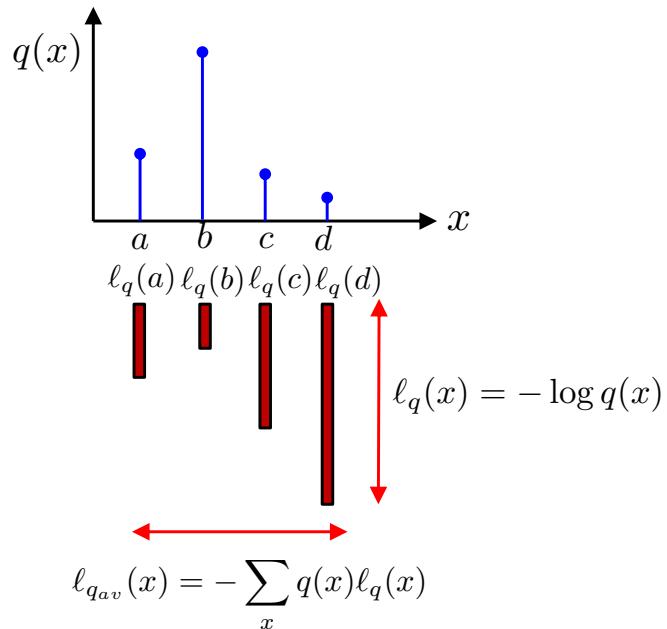
$p(x) \log_2 \frac{p(x)}{0} = +\infty, \forall p(x) > 0.$

$$D(p \parallel q) \neq D(q \parallel p)$$

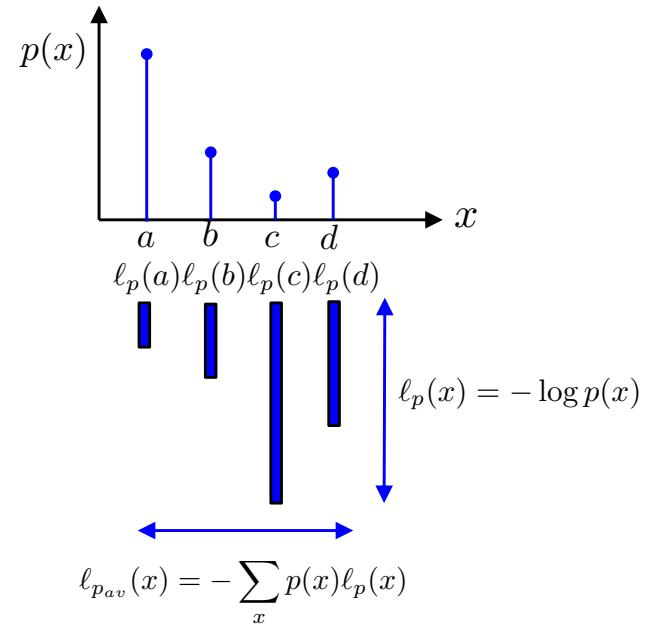
$$D(p \parallel q) \geq 0$$

Relative entropy

Assumed data distribution $q(x)$



Data distribution $p(x)$



Practical problem (the true distribution $p(x)$ is unknown)

Design code for $q(x)$

Use the designed code for $q(x)$ to code symbols from $p(x)$

Relative entropy

Under this assumption each symbol will be coded “suboptimally”

$$a \quad \text{---} \quad \ell_q(a)$$

$$\text{---} \quad \ell_p(a)$$

$$b \quad \text{---} \quad \ell_q(b)$$

$$\text{---} \quad \ell_p(b)$$

$$c \quad \text{---} \quad \ell_q(c)$$

$$\text{---} \quad \ell_p(c)$$

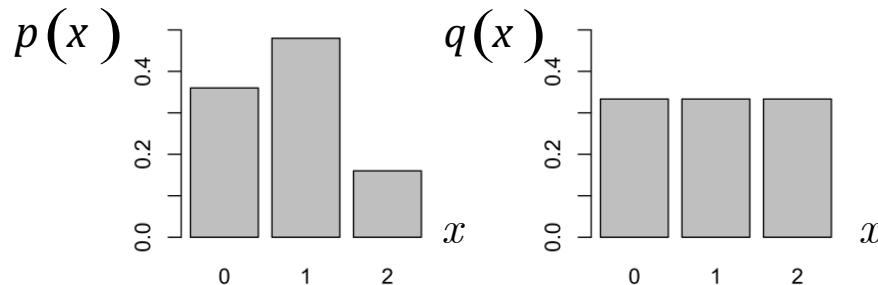
$$d \quad \text{---} \quad \ell_q(d)$$

$$\text{---} \quad \ell_p(d)$$

Question: what is the “overall loss” in this mismatched usage of codes for different sources?

Relative entropy: interpretation

- Given two distributions



- KLD “measures” in:
 - Machine learning: the level of “similarity” between these two distributions
 - Coding theory: “inefficiency” of coding designed for the pmf $q(x)$ while actually data originates from pmf $p(x)$

Real data $X \sim p(x)$

Code design $H(p(x)) + D(p(x) \parallel q(x))$
extra bits

$D(p(x) \parallel q(x)) = 0$ iff $p(x) = q(x)$

Relative entropy: interpretation

- Allow fractional code lengths $\ell_q(x) = -\log q(x)$
- Then expected length for coding $X \sim p(x)$ using $\ell_q(x)$ is

$$\begin{aligned} L &= \mathbb{E}_{X \sim p(x)} [\ell_q(X)] \\ &= - \sum_x p(x) \log q(x) \quad \xrightarrow{\text{See later as cross-entropy}} \\ &= \sum_x p(x) \log \left[\frac{p(x)}{q(x)} \frac{1}{p(x)} \right] \\ &= \sum_x p(x) \log \frac{p(x)}{q(x)} + \sum_x p(x) \log \frac{1}{p(x)} \\ &= D(p(x) \| q(x)) + H(p(x)) \end{aligned}$$

Relative entropy: asymmetry

- Example:

$$x \in \mathcal{X} = \{0,1\}; p(0) = a; p(1) = 1 - a; q(0) = b; q(1) = 1 - b.$$

$$D(p||q) = \sum_{x \in \mathcal{X}} p(x) \log_2 \frac{p(x)}{q(x)} = a \log_2 \frac{a}{b} + (1-a) \log_2 \frac{1-a}{1-b};$$

$$D(q||p) = \sum_{x \in \mathcal{X}} q(x) \log_2 \frac{q(x)}{p(x)} = b \log_2 \frac{b}{a} + (1-b) \log_2 \frac{1-b}{1-a}.$$

$$a = \frac{1}{4}; b = \frac{1}{8}, D(p||q) = \frac{1}{4} \log_2 \frac{8}{4} + \left(1 - \frac{1}{4}\right) \log_2 \left(\frac{1 - \frac{1}{4}}{1 - \frac{1}{8}} \right) = 0.0832 \text{ bit};$$

$$D(p||q) \neq D(q||p)$$

$$D(q||p) = \frac{1}{8} \log_2 \frac{4}{8} + \left(1 - \frac{1}{8}\right) \log_2 \left(\frac{1 - \frac{1}{8}}{1 - \frac{1}{4}} \right) = 0.0696 \text{ bit.}$$

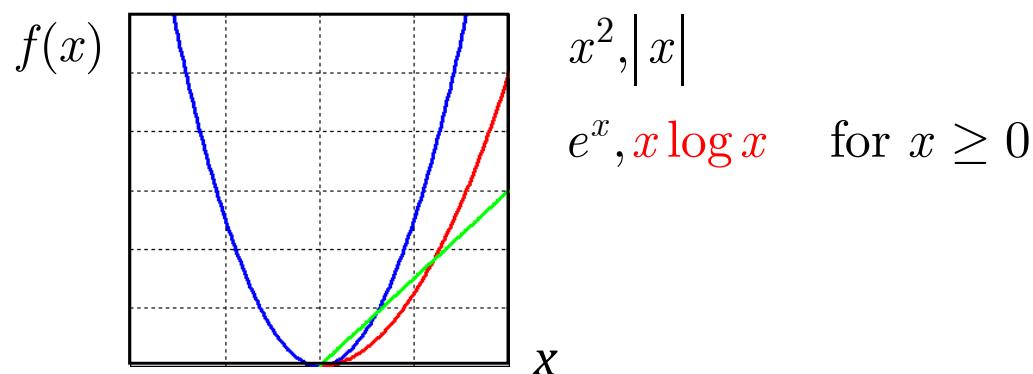
Relative entropy: recall of convexity

□

Definition (**convex function**): a function is convex (strictly convex) on the interval (a, b) , if for every $x_1, x_2 \in (a, b)$ and $0 \leq \lambda \leq 1$

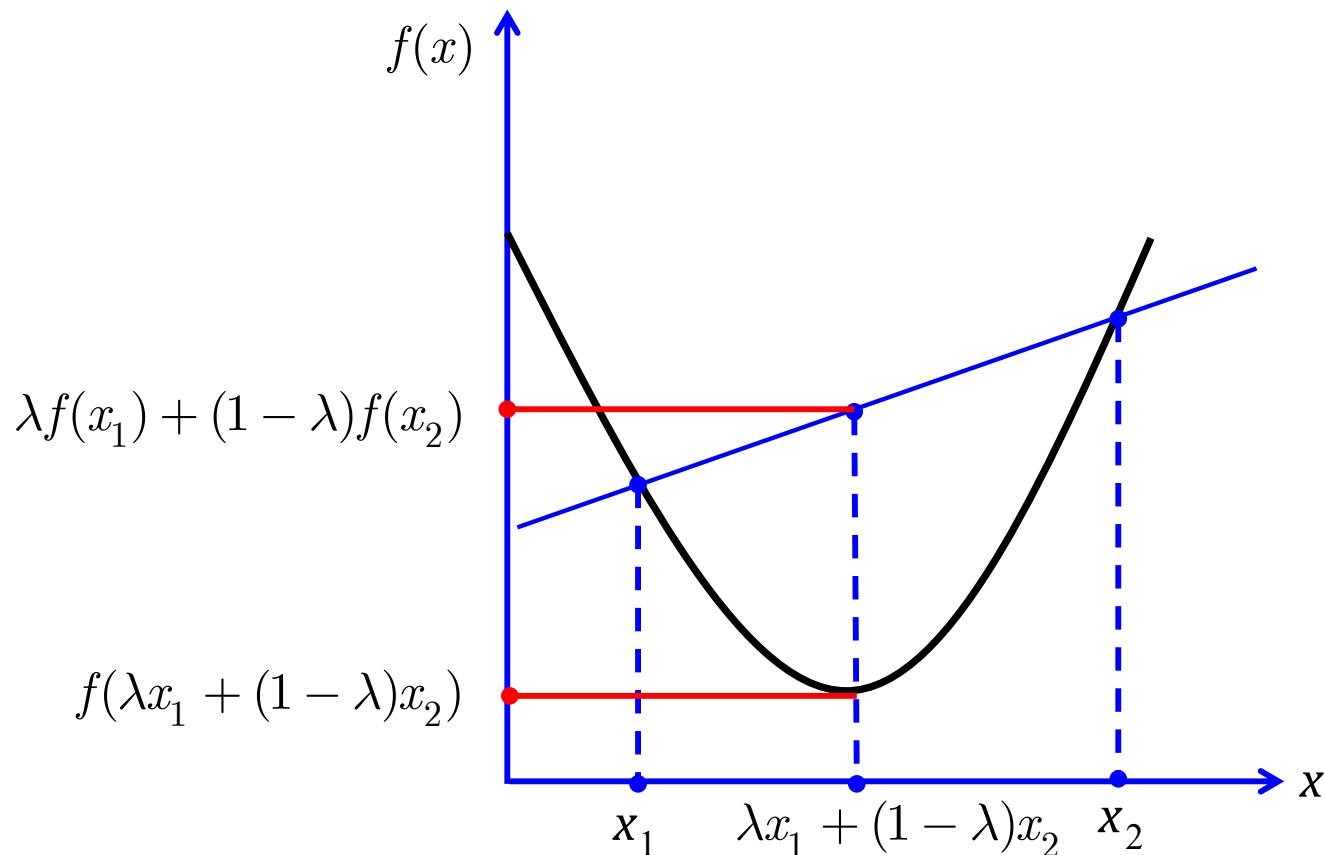
$$f(\lambda x_1 + (1 - \lambda)x_2) \leq \lambda f(x_1) + (1 - \lambda)f(x_2)$$

- If a function $f(x)$ has its second derivative to be non-negative (strictly positive), this function is convex (strictly convex).



Relative entropy: recall of convexity

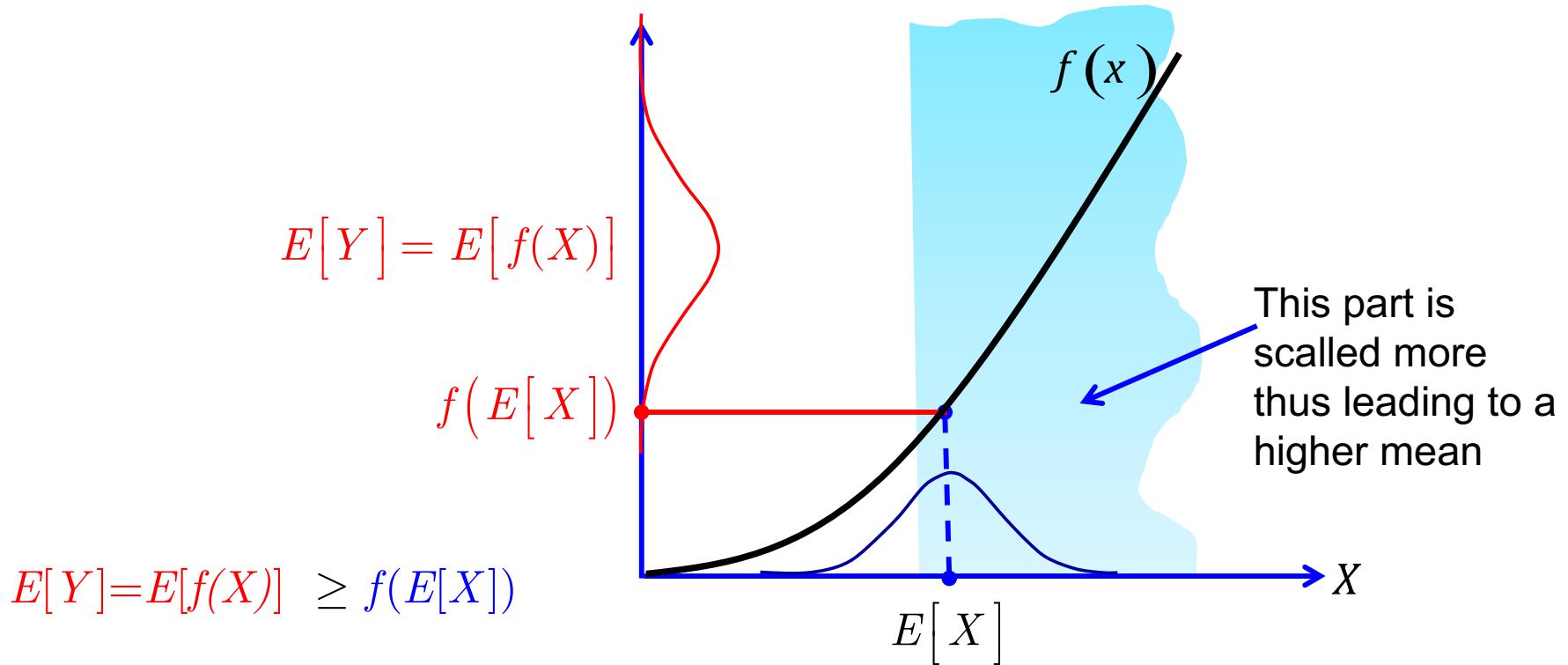
Recall: convex functions



Relative entropy: Jensen's inequality

Definition (**Jensen's inequality**): for a convex function $f(x)$ and a random variable X

$$E[f(X)] \geq f(E[X])$$



Relative entropy: log-sum inequality

□
Definition (**Log-sum inequality**): For nonnegative numbers a_1, \dots, a_N and b_1, \dots, b_N

$$\sum_{i=1}^N a_i \log\left(\frac{a_i}{b_i}\right) \geq \left(\sum_{i=1}^N a_i\right) \log \frac{\sum_{i=1}^N a_i}{\sum_{i=1}^N b_i}$$

with equality if and only if $\left(\frac{a_i}{b_i}\right) = const$

Proof:

- Assume that $a_i > 0$ and $b_i > 0$
- The function $f(t) = t \log t$ is strictly convex, since $f''(t) = \frac{1}{t} \log e > 0, t > 0$.

Relative entropy: log-sum inequality

- According to Jensen's inequality, we have:

$$\text{for } \alpha_i \geq 0, \sum_{i=1}^N \alpha_i = 1 \quad \sum_{i=1}^N \alpha_i f(t_i) \geq f\left(\sum_{i=1}^N \alpha_i t_i\right)$$

- Assume $\alpha_i = \frac{b_i}{\sum_{j=1}^N b_j}$ and $t_i = \frac{a_i}{b_i}$

$$\sum_{i=1}^N \alpha_i f(t_i) = \sum_{i=1}^N \frac{\cancel{b_i}}{\sum_{j=1}^N \cancel{b_j}} \frac{a_i}{b_i} \log\left(\frac{a_i}{b_i}\right) = \sum_{i=1}^N \frac{a_i}{\sum_{j=1}^N b_j} \log\left(\frac{a_i}{b_i}\right)$$

Relative entropy: log-sum inequality

Proof: log-sum inequality

$$\begin{aligned} & \geq f\left(\sum_{i=1}^N \alpha_i t_i\right) = \sum_{i=1}^N \alpha_i t_i \log \left(\sum_{i=1}^N \alpha_i t_i\right) = \sum_{i=1}^N \frac{b_i}{\sum_{j=1}^N b_j} \frac{a_i}{b_i} \log \left(\sum_{i=1}^N \frac{b_i}{\sum_{j=1}^N b_j} \frac{a_i}{b_i}\right) \\ & = \sum_{i=1}^N \frac{a_i}{\sum_{j=1}^N b_j} \log \left(\sum_{i=1}^N \frac{a_i}{\sum_{j=1}^N b_j}\right) \\ \Rightarrow & \sum_{i=1}^N \frac{a_i}{\sum_{j=1}^N b_j} \log \left(\frac{a_i}{b_i}\right) \geq \sum_{i=1}^N \frac{a_i}{\sum_{j=1}^N b_j} \log \left(\sum_{i=1}^N \frac{a_i}{\sum_{j=1}^N b_j}\right) \Rightarrow \sum_{i=1}^N a_i \log \left(\frac{a_i}{b_i}\right) \geq \sum_{i=1}^N a_i \log \left(\frac{\sum_{i=1}^N a_i}{\sum_{i=1}^N b_i}\right) \end{aligned}$$

Summary: properties of relative entropy

- Property 1: non-negativity of relative entropy
- Property 2: conditional relative entropy and chain rule
- Property 3: relative entropy wrt uniform pdf

Summary: properties of relative entropy

- Property 1: non-negativity of relative entropy

$$D_{\text{KL}}(p \parallel q) \geq 0$$

" = ", iff $p(x) = q(x), \forall x \in \mathcal{X}$

- Proof: $S = \{x : p(x) > 0\}$

$$\begin{aligned} -D_{\text{KL}}(p \parallel q) &= -\sum_{x \in S} p(x) \log \left(\frac{p(x)}{q(x)} \right) = \sum_{x \in S} p(x) \log \left(\frac{q(x)}{p(x)} \right) \leq \\ &\leq \left(\sum_{x \in S} p(x) \right) \log \frac{\sum_{x \in S} (q(x))}{\sum_{x \in S} (p(x))} = 1 \cdot \log \frac{1}{1} = 0 \quad \sum_i a_i \log \left(\frac{a_i}{b_i} \right) \geq \left(\sum_i a_i \right) \log \frac{\sum_i a_i}{\sum_i b_i} \quad \text{Log-sum INQ} \\ &\quad D_{\text{KL}}(p \parallel q) \geq 0 \quad D_{\text{KL}}(p \parallel q) = 0 \\ &\quad p(x) = q(x) \end{aligned}$$

$$-D(p(x) \parallel q(x)) = \sum p(x) \log \left(\frac{q(x)}{p(x)} \right) \leq \log \left(\sum p(x) \frac{q(x)}{p(x)} \right) = \log \left(\sum q(x) \right) = 0$$

Summary: properties of relative entropy

- Property 2: conditional relative entropy and chain rule

- Conditional relative entropy

$$D_{\text{KL}}(p(y|x) \| q(y|x)) = \sum_{x \in \mathcal{X}} p(x) \sum_{y \in \mathcal{Y}} p(y|x) \log_2 \frac{p(y|x)}{q(y|x)} = E_{p(x,y)} \left[\log_2 \frac{p(y|x)}{q(y|x)} \right]$$

- Chain rule for relative entropy

$$D_{\text{KL}}(p(x,y) \| q(x,y)) = D_{\text{KL}}(p(x) \| q(x)) + D_{\text{KL}}(p(y|x) \| q(y|x))$$

Summary: properties of relative entropy

- Proof:

$$\begin{aligned} D_{\text{KL}}(p(x,y) \parallel q(x,y)) &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x,y) \log_2 \frac{p(x,y)}{q(x,y)} \\ &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x,y) \log_2 \left[\frac{p(x)}{q(x)} \frac{p(y|x)}{q(y|x)} \right] && \Bigg| \quad p(x,y) = p(x)p(y|x) \\ &= \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x,y) \log_2 \frac{p(x)}{q(x)} + \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x,y) \log_2 \frac{p(y|x)}{q(y|x)} && \Bigg| \quad \log ab = \log a + \log b \\ &= D_{\text{KL}}(p(x) \parallel q(x)) + D_{\text{KL}}(p(y|x) \parallel q(y|x)) \quad \square \end{aligned}$$

Summary: properties of relative entropy

- Property 3: relative entropy wrt uniform pmf

- Assume a random variable with the **uniform distribution** ($q(x) = 1/|\mathcal{X}|$)
- Then, KLD between any $p(x)$ and uniform $q(x)$ is:

$$D_{\text{KL}}(p\|q) = \log_2 |\mathcal{X}| - H(X)$$

$$D_{\text{KL}}(p\|q) = \sum_{x \in \mathcal{X}} p(x) \log_2 \frac{p(x)}{q(x)} = \underbrace{\sum_{x \in \mathcal{X}} p(x) \log_2 p(x)}_{-H(X)} - \underbrace{\sum_{x \in \mathcal{X}} p(x) \log_2 q(x)}_{1} \underbrace{1}_{1/|\mathcal{X}|}$$

$$\Rightarrow H(X) = \log_2 |\mathcal{X}| - D_{\text{KL}}(p\|q)$$

- Thus, the distribution $p(x)$ that has the maximum entropy $H(X) = \log_2 |\mathcal{X}|$ is uniform, i.e., since $D_{\text{KL}}(p\|q) \geq 0$

Cross-entropy

□
Definition (**cross-entropy**): *Cross-entropy* between pmfs $p(x)$ and $q(x)$:

$$H(p, q) = - \sum_{x \in \mathcal{X}} p(x) \log q(x) = -E_{p(x)}[\log q(x)]$$

- A link to KLD:

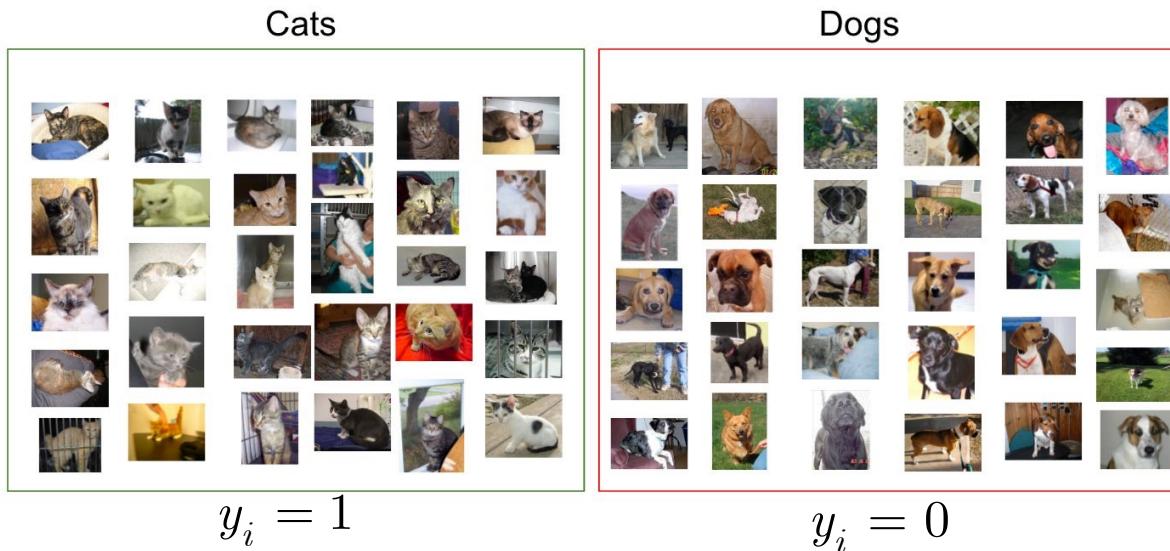
$$H(p, q) = H(p) + D_{\text{KL}}(p \parallel q)$$

- Proof:

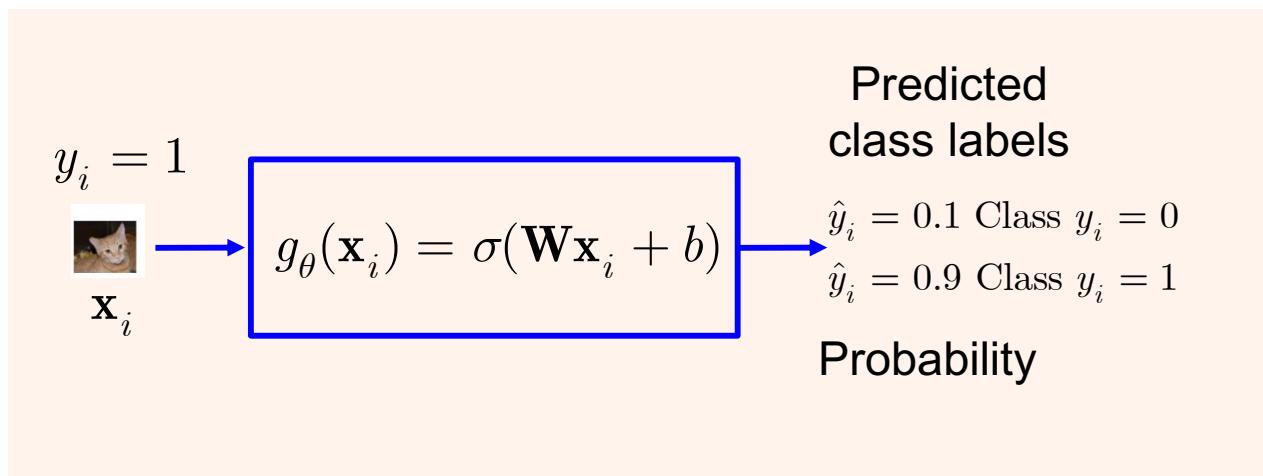
$$\begin{aligned} D_{\text{KL}}(p \parallel q) &= E_{p(x)} \left[\log \frac{p(x)}{q(x)} \right] = E_{p(x)} [\log p(x)] - E_{p(x)} [\log q(x)] \\ &= -H(p) + H(p; q) \end{aligned}$$

- Remark: since $H(p)$ is constant, optimizing $H(p, q)$ wrt $q(x)$ is equivalent to $D_{\text{KL}}(p \parallel q)$

Cross-entropy: in binary classification



- Total N training samples
- All samples are labeled $\{\mathbf{x}_i, y_i\}_{i=1}^N$
- We have two labels $y_i \in \{0, 1\}$



- Goal:
Train the classifier $g_{\theta}(\mathbf{x}_i)$ producing the most accurate class prediction

Cross-entropy: in binary classification

- Notations:
 - True label $y_i \in \{0,1\}$
 - Predicted label $\hat{y}_i = p_\theta(y_i = 1 | \mathbf{x}_i) = g_\theta(\mathbf{x}_i)$
 $1 - \hat{y}_i = p_\theta(y_i = 0 | \mathbf{x}_i) = 1 - g_\theta(\mathbf{x}_i)$
- Log-likelihood for N independent training samples:

$$\begin{aligned}\ell(\theta) &= \log L(\theta) = \log \prod_{i=1}^N \underbrace{\left[g_\theta(\mathbf{x}_i) \right]}_{\hat{y}_i}^{(y_i)} \underbrace{\left[1 - g_\theta(\mathbf{x}_i) \right]}_{1-\hat{y}_i}^{(1-y_i)} \\ &= \sum_{i=1}^N y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)\end{aligned}$$

- Optimal estimation of parameter according to Maximum Likelihood (ML)

$$\hat{\theta} = \arg \max_{\theta} \ell(\theta) = \arg \min_{\theta} -\ell(\theta) = \arg \min_{\theta} \sum_{i=1}^N \underbrace{-(y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i))}_{H(y_i; \hat{y}_i)}$$

Cross-entropy: in binary classification

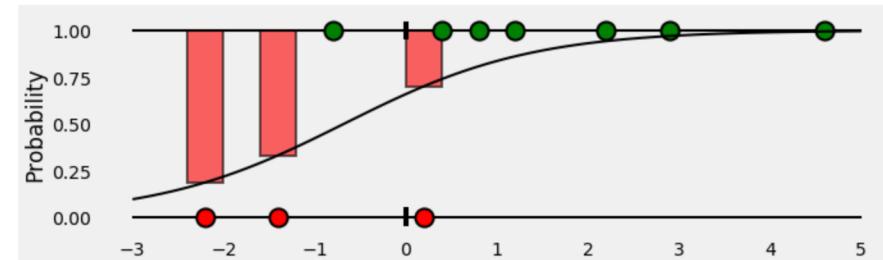
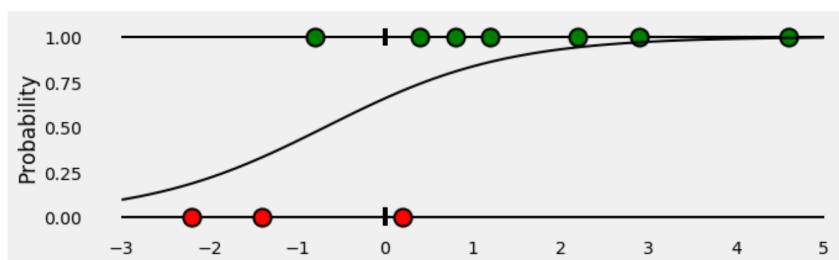
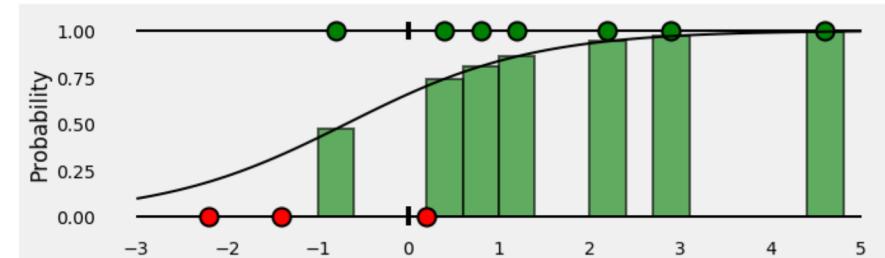
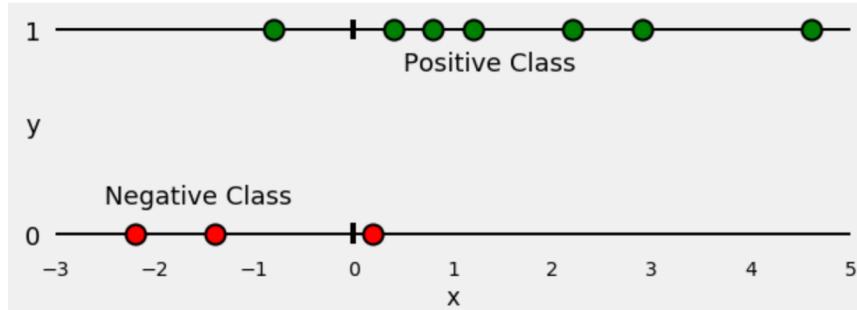
- Summary:
 - True label pmf $p \in \{y, 1 - y\}$
 - Predicted label pmf $q \in \{\hat{y}, 1 - \hat{y}\}$

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N H(y_i, \hat{y}_i) = -\frac{1}{N} \sum_{i=1}^N [y_i \log \hat{y}_i + (1 - y_i) \log (1 - \hat{y}_i)]$$

- Link to KLD (entropy is constant and it is equivalent)

$$H(p, q) = \underbrace{H(p)}_{const} + D_{KL}(p \parallel q)$$

Cross-entropy: in binary classification



<https://towardsdatascience.com/understanding-binary-cross-entropy-log-loss-a-visual-explanation-a3ac6025181a>

Mutual information

Definition (mutual information): *Mutual information* between two r.v. X and Y is defined by:

$$I(X;Y) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{X,Y}(x,y) \log \frac{p_{X,Y}(x,y)}{p_X(x)p_Y(y)} = E_{p_{X,Y}(x,y)} \left[\log \frac{p_{X,Y}(x,y)}{p_X(x)p_Y(y)} \right]$$
$$I(X;Y) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p_{X,Y}(x,y) \log \frac{p_{X|Y}(x|y)}{p_X(x)} \cancel{\frac{p_Y(y)}{p_Y(y)}}$$

$$I(X;Y) = I(Y;X)$$

$$I(X;Y) \geq 0$$

$$\Rightarrow I(X;Y) = 0 \Leftrightarrow X \perp Y \quad (\text{if independent})$$

Mutual information

Mutual information via KL-divergence

$$I(X;Y) = D_{\text{KL}}(p_{X,Y}(x,y) \parallel p_X(x)p_Y(y)) = E_{p_{X,Y}(x,y)} \left[\log \frac{p_{X,Y}(x,y)}{p_X(x)p_Y(y)} \right]$$

$$I(X;Y) = \mathbb{E}_{p(y)} [D_{\text{KL}}(p(x|y) \parallel p(x))] = \mathbb{E}_{p(x)} [D_{\text{KL}}(p(y|x) \parallel p(y))]$$

- Mutual information measures how different the joint pmf is from the product of marginal pmfs, i.e., a **real pmf**, vs an assumption of **independent variables**

Pay attention: on the equivalence of notations $p_{X,Y}(x,y) = p(x,y)$

Mutual information

Mutual information via KL-divergence

$$\begin{aligned} I(X;Y) &= D_{\text{KL}}(p(x,y) \| p(x)p(y)) = \mathbb{E}_{p(x,y)} \left[\log \frac{p(x,y)}{p(x)p(y)} \right] \\ &= \mathbb{E}_{p(x)p(y|x)} \left[\log \frac{p(y|x)p(x)}{p(x)p(y)} \right] = \mathbb{E}_{p(x)} \left[\mathbb{E}_{p(y|x)} \left[\log \frac{p(y|x)}{p(y)} \right] \right] \\ &= \mathbb{E}_{p(x)} [D_{\text{KL}}(p(y|x) \| p(y))] \end{aligned}$$

Summary: properties of mutual information

- Property 1: non-negativity of mutual information
- Property 2: mutual information decomposition via entropy
- Property 3: conditional mutual information
- Property 4: chain rule for mutual information
- Property 5: non-creativity of information by processing /data processing inequality/

Mutual information: properties

- Property 1: non-negativity of mutual information (from non-negativity of KLD)

$$I(X;Y) \geq 0, \quad \Rightarrow I(X;Y) = 0 \Leftrightarrow X \perp Y$$

$$I(X;Y,Z) \geq 0 \quad \text{If } X \text{ and } Y \text{ are independent}$$

$$I(X;Y) = H(X) - H(X|Y) \quad H(X) \geq 0, H(X|Y) \geq 0$$

$$H(X) \geq H(X|Y) \Rightarrow I(X;Y) \geq 0$$

- Property 2: mutual information decomposition via entropy

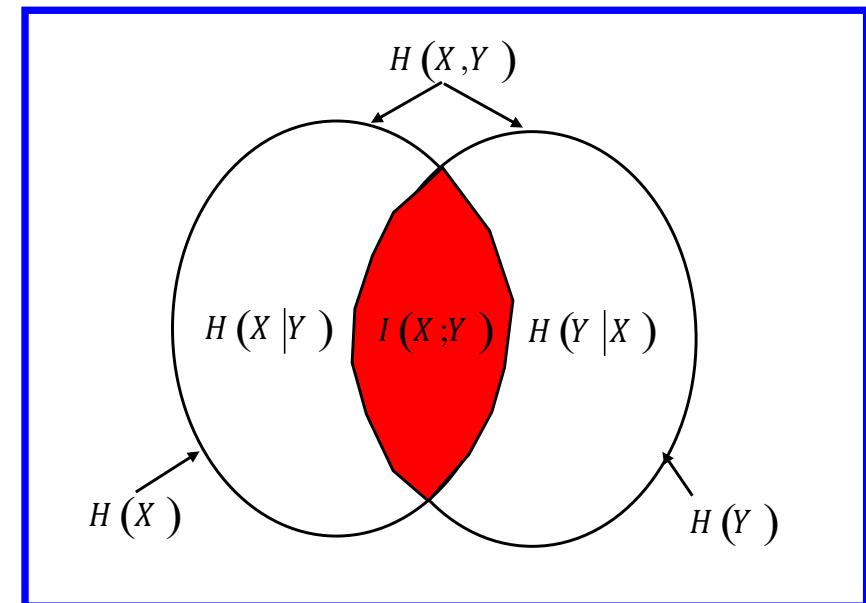
$$I(X;Y) = H(X) - H(X|Y)$$

$$I(X;Y) = H(Y) - H(Y|X)$$

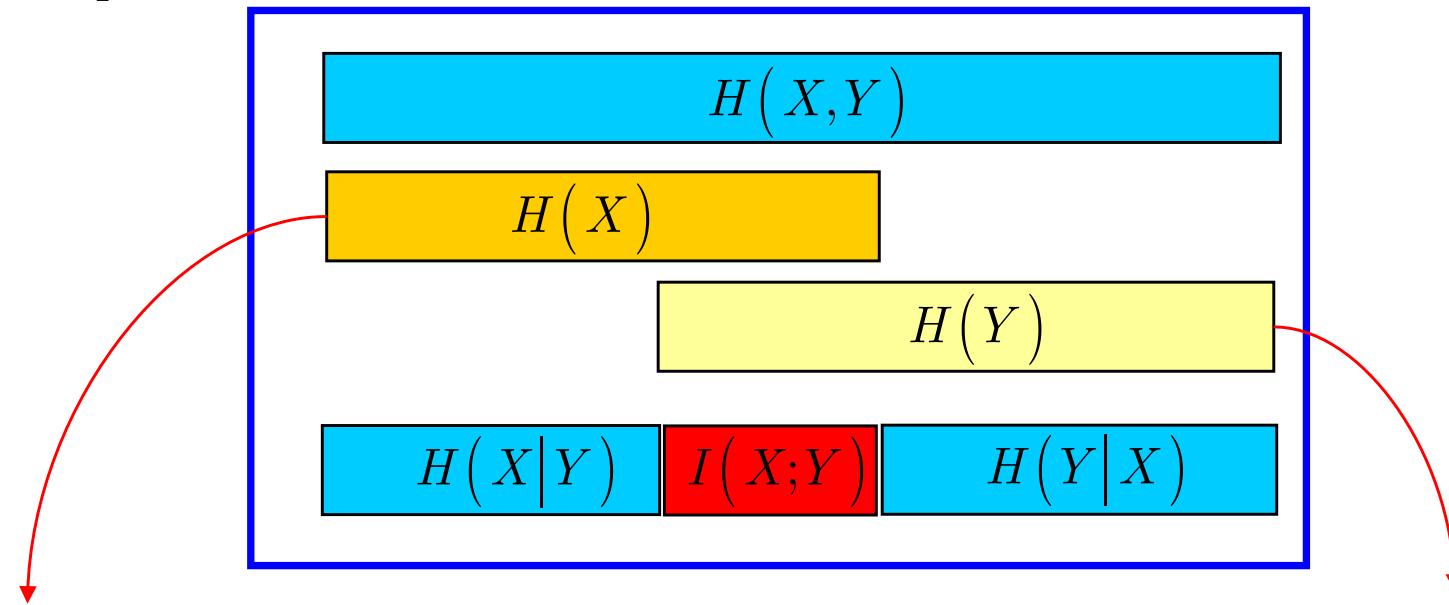
$$I(X;Y) = H(X) + H(Y) - H(X,Y)$$

$$I(X;X) = H(X)$$

$$I(X;X) = H(X) - H(X|X) \underbrace{}_0$$



Mutual information: properties



$$I(X;Y) = H(X) - H(X|Y) \quad I(X;Y) = H(Y) - H(Y|X)$$

Démo:

$$I(X;Y) = E_{p_{X,Y}} \left[\log \frac{p_{X|Y}(x|y)}{p_X(x)} \right] = \underbrace{E_{p_{X,Y}} \left[\log p_{X|Y}(x|y) \right]}_{-H(X|Y)} - \underbrace{E_{p_{X,Y}} \left[\log p_X(x) \right]}_{H(X)}$$

Summary: properties of mutual information

- Property 3: conditional mutual information

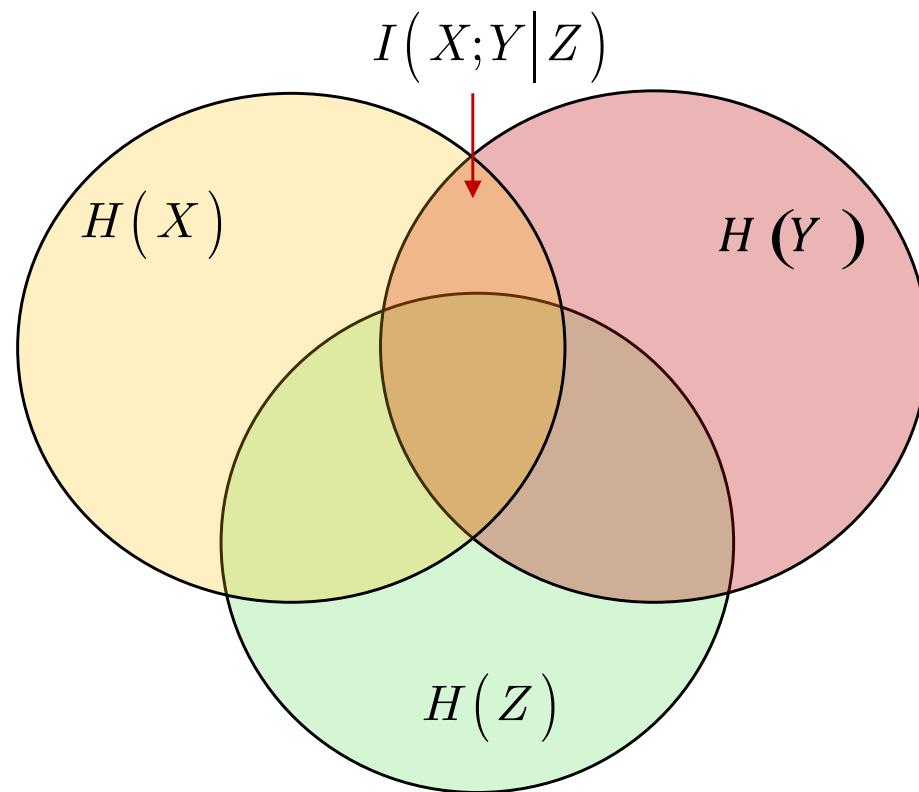
$$I(X;Y|Z) = H(X|Z) - H(X|Y,Z)$$

$$\begin{aligned} I(X;Y|Z) &= E_{p(x,y,z)} \left[\log \frac{p_{X,Y|Z}(x,y|z)}{p_{X|Z}(x|z)p_{Y|Z}(y|z)} \right] = E_{p(x,y,z)} \left[\log \frac{p_{X|Y,Z}(x|y,z)}{p_{X|Z}(x|z)} \right] \\ &\quad p_{X,Y|Z}(x,y|z) = p_{Y|Z}(y|z)p_{X|Y,Z}(x|y,z) \\ &= \underbrace{E_{p(x,y,z)} \left[\log p_{X|Y,Z}(x|y,z) \right]}_{-H(X|Y,Z)} - \underbrace{E_{p(x,y,z)} \left[\log p_{X|Z}(x|z) \right]}_{H(X|Z)} \end{aligned}$$

Summary: properties of mutual information

- Property 3: conditional mutual information

$$I(X;Y|Z) = H(X|Z) - H(X|Y,Z)$$



Summary: properties of mutual information

- Property 4: chain rule for mutual information

$$I(X_1, X_2, \dots, X_N; Y) = \sum_{i=1}^N I(X_i; Y | X_{i-1}, \dots, X_1)$$

$$I(X_1, X_2, \dots, X_N; Y) = H(X_1, X_2, \dots, X_N) - H(X_1, X_2, \dots, X_N | Y)$$

$$\left| \begin{array}{l} I(Z; Y) = H(Z) - H(Z | Y) \\ \text{Chain rule} \\ H(X_1, X_2, \dots, X_N) = \sum_{i=1}^N H(X_i | X_{i-1}, \dots, X_1) \end{array} \right.$$

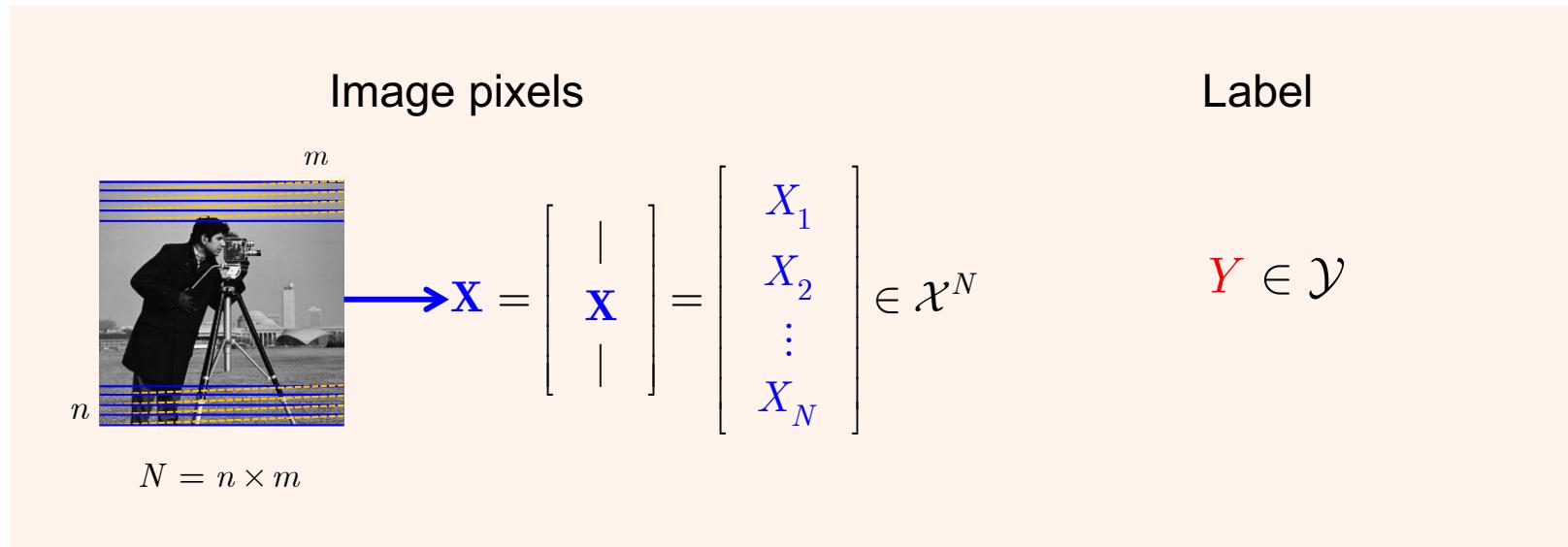
$$= \sum_{i=1}^N H(X_i | X_{i-1}, \dots, X_1) - \sum_{i=1}^N H(X_i | X_{i-1}, \dots, X_1, Y) =$$

$$= \sum_{i=1}^N I(X_i; Y | X_{i-1}, \dots, X_1)$$

Summary: properties of mutual information

- Property 4: chain rule for mutual information

$$I(X_1, X_2, \dots, X_N; Y) = \sum_{i=1}^N I(X_i; Y | X_{i-1}, \dots, X_1)$$



Mutual information: properties

- Example: more complex decomposition

$$I(X_1, X_2; Y_1, Y_2) = I(X_1; Y_1, Y_2) + I(X_2; Y_1, Y_2 | X_1)$$

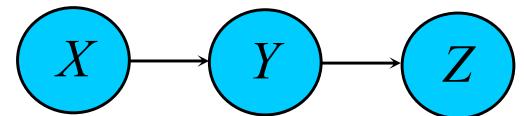
$$I(X_1; Y_1, Y_2) = I(X_1; Y_1) + I(X_1; Y_2 | Y_1) = I(X_1; Y_2) + I(X_1; Y_1 | Y_2)$$

$$\begin{aligned} I(X_2; Y_1, Y_2 | X_1) &= I(X_2; Y_1 | X_1) + I(X_2; Y_2 | X_1, Y_1) \\ &= I(X_2; Y_2 | X_1) + I(X_2; Y_1 | X_1, Y_2) \end{aligned}$$

Summary: properties of mutual information

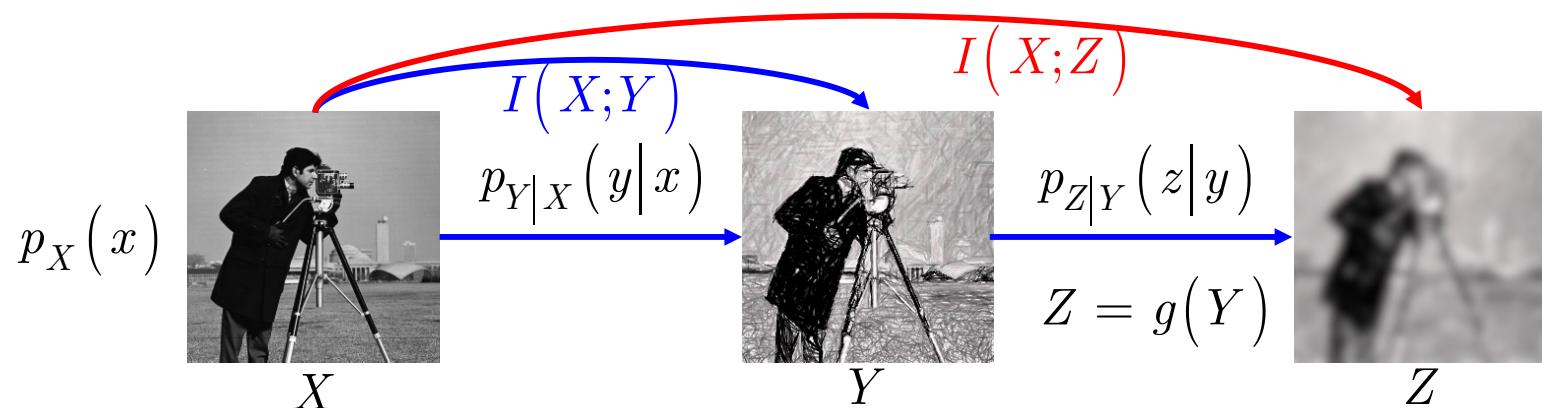
- Property 5: non-creativity of information by processing /data processing inequality/

- Consider the Markov chain $p_{X,Y,Z}(x,y,z) = p_X(x)p_{Y|X}(y|x)p_{Z|Y}(z|y)$



- Two inequalities are valid

$$I(X;Y) \geq I(X;Z) \quad I(Y;Z) \geq I(X;Z)$$



Summary: properties of mutual information

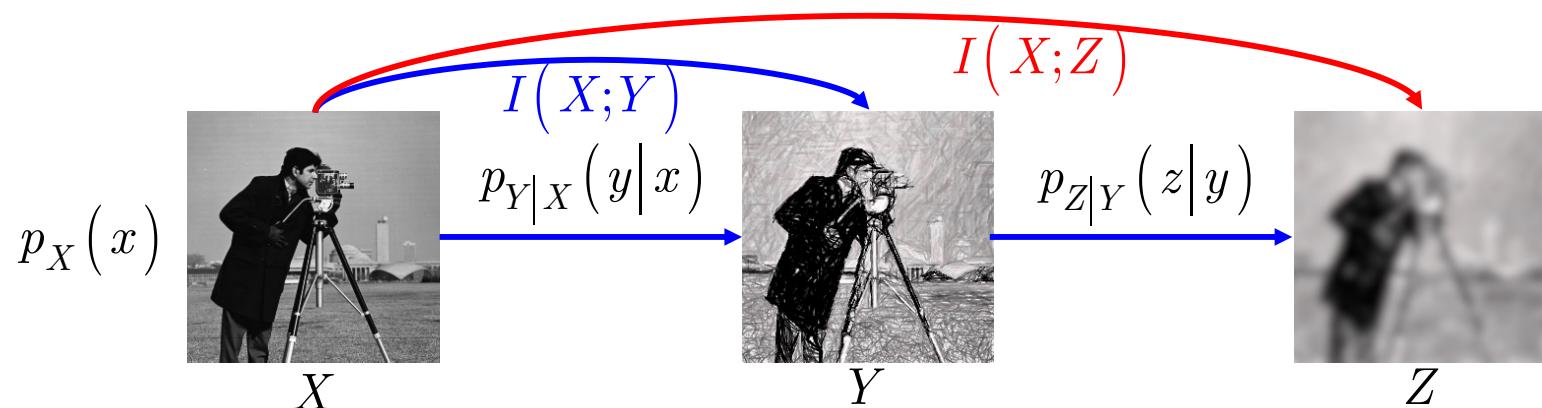
- **Data processing inequality**

The amount of information cannot be increased by any processing $Z = g(Y)$

This also includes deep networks!

- Two inequalities are valid

$$I(X;Y) \geq I(X;Z) \quad I(Y;Z) \geq I(X;Z)$$



Summary: properties of mutual information

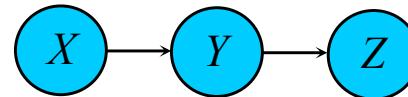
- **Proof of data processing inequality**

- Using chain rule for mutual information

$$I(X;Y,Z) = I(X;Z) + I(X;Y|Z) \quad (a)$$

$$= I(X;Y) + I(X;Z|Y) \quad (b)$$

$$I(X;Z|Y) = 0$$

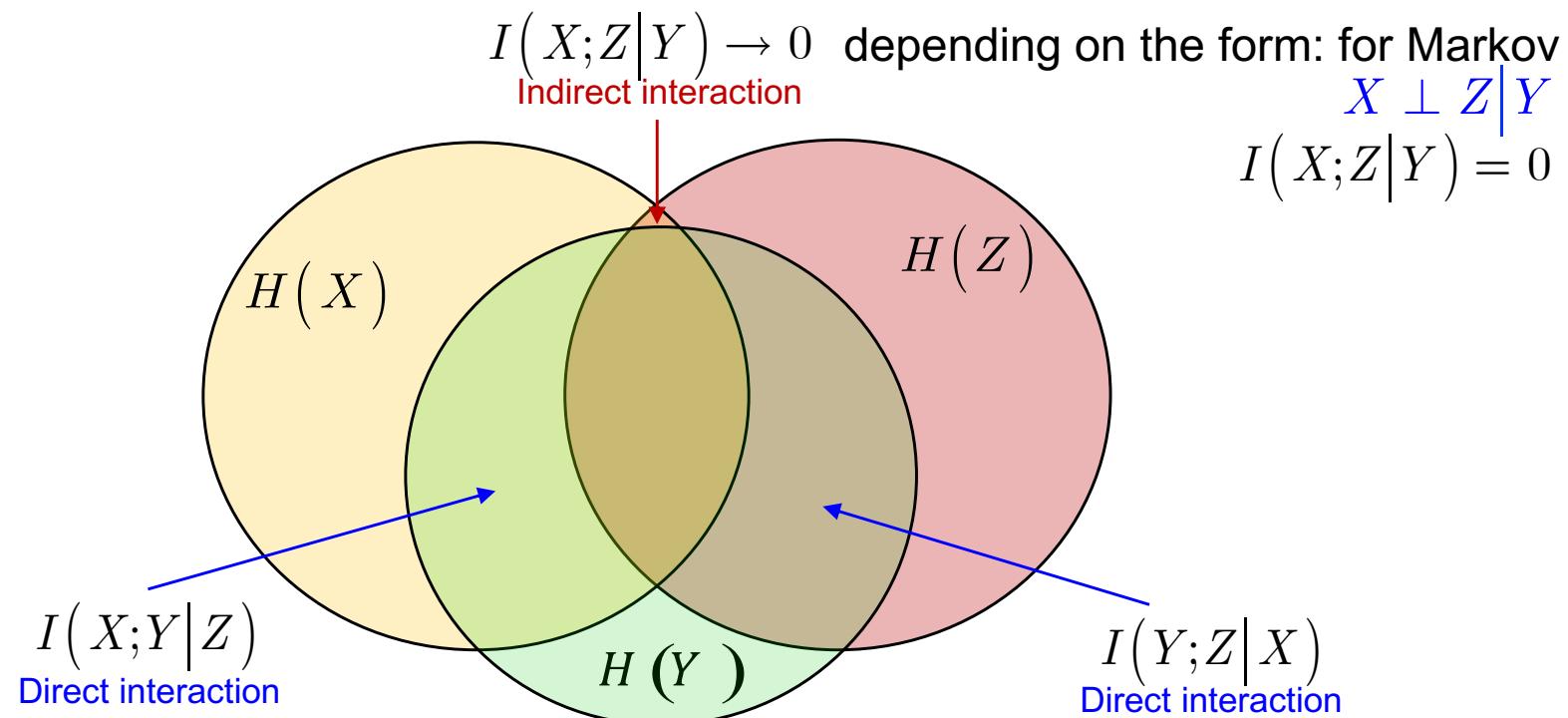
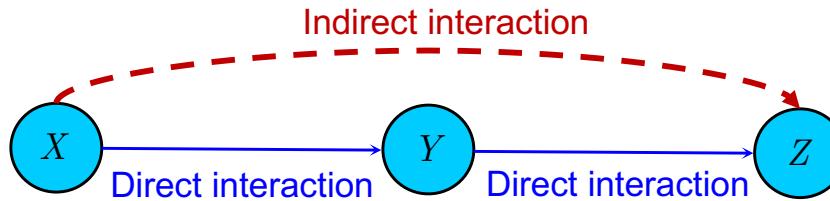


$$X \perp Z|Y$$

$$I(X;Y) = I(X;Z) + I(X;Y|Z) \Rightarrow I(X;Y) \geq I(X;Z)$$

Summary: properties of mutual information

- Consider a general case



Scope

Part II: continuous random variables and vectors

Scope

- Information theoretic measures for discrete/continuous variables

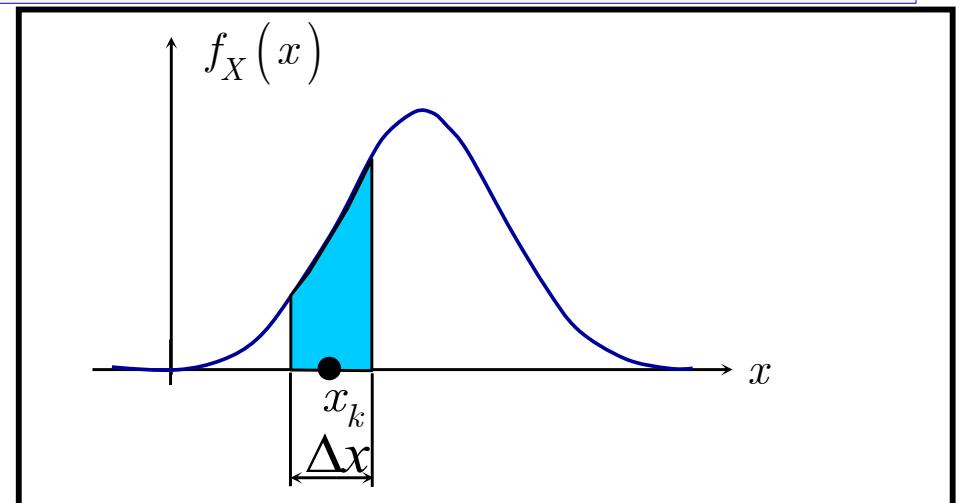
- Entropy → Differential entropy
- Conditional entropy
- Joint entropy
- Relative entropy (KL-divergence)
- Cross entropy
- Mutual information
- Additional topics:
 - f-divergence
 - variational inference
 - practical computations from samples

Continuous R.Vs.: Differential entropy

Definition (entropy): *Differential entropy* of r.v. X with pdf $f_X(x)$

$$h(X) = - \int_{\mathcal{X}} f_X(x) \log_2 f_X(x) dx = E_{f_X}[-\log_2 f_X(x)]$$

$$P[x_k; \Delta x] = \int_{\Delta x} f_X(x) dx \cong f_X(x_k) \Delta x$$



Continuous R.Vs.: Differential entropy

$$\begin{aligned} H(X; \Delta x) &= -\sum_{i=1}^N P[x_k, \Delta x] \log_2 P[x_k, \Delta x] = -\sum_{i=1}^N f_X(x_k) \Delta x \log_2 f_X(x_k) \Delta x \\ &= -\sum_{i=1}^N f_X(x_k) \Delta x \log_2 f_X(x_k) - \underbrace{\sum_{i=1}^N f_X(x_k) \Delta x \log_2 \Delta x}_{= 1 \quad const} \\ &\quad \int_{-\infty}^{+\infty} f_X(x) dx = 1 \end{aligned}$$

In the limit Δx tends to zero for large N . As a result, $\log_2 \Delta x$ tends to infinity.

$$H(X; \Delta x) = \underbrace{-\int f_X(x) \log_2 f_X(x) dx}_{h(X)} - \log_2 \Delta x.$$

Continuous R.Vs.: Differential entropy

□ **Definition (Differential entropy of Gaussian r.v.):** *differential entropy of Gaussian r.v.*

$X \sim \mathcal{N}(0, \sigma_X^2)$ is:

$$h(X) = -E_{f_X} [\log_2 f_X(x)] = \frac{1}{2} \log_2 (2\pi e \sigma_X^2)$$

- Proof:

$$\begin{aligned} h(X) &= - \int_{-\infty}^{\infty} f_X(x) \ln f_X(x) dx \quad [nats] = - \int_{-\infty}^{\infty} f_X(x) \left[-\frac{x^2}{2\sigma_X^2} - \ln \sqrt{2\pi\sigma_X^2} \right] dx = \\ &= \int_{-\infty}^{\infty} f_X(x) \frac{x^2}{2\sigma_X^2} dx + \ln \sqrt{2\pi\sigma_X^2} \int_{-\infty}^{\infty} f_X(x) dx \end{aligned}$$

Continuous R.Vs.: Differential entropy

$$h(X) = \int_{-\infty}^{\infty} f_X(x) \frac{x^2}{2\sigma_X^2} dx + \ln \sqrt{2\pi\sigma_X^2} \underbrace{\int_{-\infty}^{\infty} f_X(x) dx}_{=1}$$

$$Var[X] = \int_{-\infty}^{\infty} x^2 f_X(x) dx = \sigma_X^2$$

$$\int_{-\infty}^{\infty} f_X(x) \frac{x^2}{2\sigma_X^2} dx = \frac{\sigma_X^2}{2\sigma_X^2} = \frac{1}{2}$$

$$h(X) = \frac{1}{2} + \ln \sqrt{2\pi\sigma_X^2} = \frac{1}{2} \ln e + \frac{1}{2} \ln 2\pi\sigma_X^2 = \frac{1}{2} [\ln e + \ln 2\pi\sigma_X^2] = \frac{1}{2} \ln(2\pi e \sigma_X^2) [nats]$$

$$h(X) = -E_{f_X} [\log_2 f_X(x)] = \frac{1}{2} \log_2 (2\pi e \sigma_X^2) [\text{bits}]$$

Continuous R.Vs.: Differential entropy

- Differential entropy can be **negative**

$$h(X) = \frac{1}{2} \log_2(2\pi e \sigma_X^2)$$

$$\log_2(a) \leq 0, \text{ for } a \leq 1$$

$$h(X) \leq 0, \text{ for } 2\pi e \sigma_X^2 \leq 1 \Rightarrow \sigma_X^2 \leq \frac{1}{2\pi e}$$



- Entropy of discrete random variables $H(X) \geq 0$

Continuous R.Vs.: Differential entropy of sum

Definition (Differential entropy of sum of two independent Gaussian r.v.s): given two independent random variables $X_1 \sim \mathcal{N}\left(\bar{x}_1, \sigma_{X_1}^2\right)$ and $X_2 \sim \mathcal{N}\left(\bar{x}_2, \sigma_{X_2}^2\right)$, their sum $Y = X_1 + X_2$ has the differential entropy

$$h(Y) = E_{f_Y}[-\log_2 f_Y(y)] = \frac{1}{2} \log_2 \left(2\pi e \left(\sigma_{X_1}^2 + \sigma_{X_2}^2 \right) \right)$$

Continuous R.Vs.: Differential entropy

□
Definition (Differential entropy of uniform r.v.): *differential entropy of uniform r.v.*

X defined on $[a, b]$ is:

$$h(X) = \log_2(b - a)$$

- Proof:

$$h(X) = - \int_{\mathcal{X}} f_X(x) \log_2 f_X(x) dx = - \int_a^b \frac{1}{b-a} \log_2 \frac{1}{b-a} dx = \log_2(b-a)$$

Continuous R.Vs.: Differential entropy

□

Definition (Differential entropy of exponential r.v.): *differential entropy of exponential r.v.* X is:

$$h(X) = \log_2 \frac{e}{\lambda}$$

▪ **Proof:**

$$h(X) = - \int_x f_X(x) \log_2 f_X(x) dx = - \int_0^\infty \lambda \exp(-\lambda x) \log_2(\lambda \exp(-\lambda x)) dx =$$

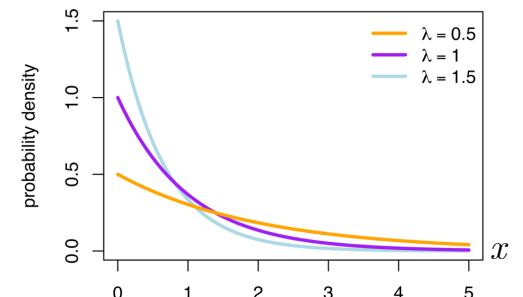
$$= - \int_0^\infty \lambda \exp(-\lambda x) (\log_2(\lambda) - \lambda x \log_2(e)) dx = I_1 + I_2;$$

$$I_1 = -\log_2(\lambda) \underbrace{\int_0^\infty \lambda \exp(-\lambda x) dx}_{=1} = \log_2 \frac{1}{\lambda};$$

$$I_2 = - \int_0^\infty -(\lambda x \log_2(e)) \lambda \exp(-\lambda x) dx = \lambda \log_2(e) \underbrace{\int_0^\infty x \lambda \exp(-\lambda x) dx}_{E[X]=1/\lambda} =$$

$$= \lambda \log_2(e) E[X] = \log_2(e);$$

$$I_1 + I_2 = \log_2 \frac{1}{\lambda} + \log_2(e) = \log_2 \frac{e}{\lambda} = 1 - \ln \lambda$$



Continuous R.Vs.: Conditional differential entropy

□
Definition (Conditional differential entropy): *conditional differential entropy* is defined as

$$h(X|Y) = - \int_{\mathcal{X}} \int_{\mathcal{Y}} f_{X,Y}(x,y) \log_2 f_{X|Y}(x|y) dx dy = h(X,Y) - h(Y)$$

$$f_{X|Y}(x|y) = f_{X,Y}(x,y) / f_Y(y)$$

$$h(X|Y) \neq h(Y|X)$$

Continuous R.Vs.: Joint differential entropy

Definition (Joint differential entropy): *joint differential entropy* is defined as

$$\begin{aligned} h(X, Y) &= - \int_{\mathcal{X}} \int_{\mathcal{Y}} f_{X,Y}(x, y) \log_2 f_{X,Y}(x, y) dx dy = \\ &= h(X) + h(Y|X) = h(Y) + h(X|Y) \end{aligned}$$

$$h(X, Y) = h(Y, X)$$

Continuous R.Vs.: max entropy for fixed variance

- Gaussian has maximum differential entropy **among all distributions with the same variance**
- Let $q(x)$ be any density satisfying $\int x^2 q(x) dx = \sigma^2$. Let $p(x) = \mathcal{N}(0, \sigma_X^2)$. Then
$$h(q(x)) \leq h(p(x))$$
- We will show two possible approaches to prove this:
 - Approach 1 based on optimization
 - Approach 2 based on KLD properties
- Remark: for simplicity we will use distributions $p(x)$ and $q(x)$ to denote pdfs.

Continuous R.Vs.: max entropy (app 1)

- Recall (we have used it for the uniform distribution)
- We will use the method of Lagrange multipliers (unconstrained regularization)

$$\mathcal{L}(x_1, x_2, \dots, x_N) = f(x_1, x_2, \dots, x_N) + \lambda(g(x_1, x_2, \dots, x_N) - k)$$

$$\left(\hat{x}_1, \hat{x}_2, \dots, \hat{x}_N \right) = \arg \min_{x_1, x_2, \dots, x_N} \mathcal{L}(x_1, x_2, \dots, x_N)$$

- Solution $\frac{\partial \mathcal{L}(x_1, x_2, \dots, x_N)}{\partial x_i} = \frac{\partial f(x_1, \dots, x_i, \dots, x_N)}{\partial x_i} + \lambda \frac{\partial g(x_1, \dots, x_i, \dots, x_N)}{\partial x_i} = 0$

$$\frac{d\mathcal{L}(x_1, x_2, \dots, x_N)}{d\lambda} = g(x_1, x_2, \dots, x_N) - k = 0$$

Continuous R.Vs.: max entropy (app 1)

- We should satisfy the constraints:

$$\int_{-\infty}^{+\infty} q(x) dx = 1 \quad (a)$$

$$\int_{-\infty}^{+\infty} (x - \mu)^2 q(x) dx = \sigma^2 \quad (b)$$

- The Lagrangian is:

$$\tilde{\mathcal{L}}(q(x)) = - \int_{-\infty}^{+\infty} q(x) \ln q(x) dx + \lambda_0 \left(1 - \int_{-\infty}^{+\infty} q(x) dx \right) + \lambda_1 \left(\sigma^2 - \int_{-\infty}^{+\infty} (x - \mu)^2 q(x) dx \right)$$

- The maximization of $\tilde{\mathcal{L}}(q(x))$ is equivalent to minimization of $\mathcal{L}(q(x)) = (-1) \times \tilde{\mathcal{L}}(q(x))$

$$\mathcal{L}(q(x)) = \int_{-\infty}^{+\infty} q(x) \ln q(x) dx - \lambda_0 \left(1 - \int_{-\infty}^{+\infty} q(x) dx \right) - \lambda_1 \left(\sigma^2 - \int_{-\infty}^{+\infty} (x - \mu)^2 q(x) dx \right)$$

Continuous R.Vs.: max entropy (app 1)

- The derivative wrt $q(x)$ and setting it to zero:

$$\frac{d\mathcal{L}(q(x))}{dq(x)} = 1 + \ln q(x) + \lambda_0 + \lambda_1(x - \mu)^2 = 0$$

$$\ln q(x) = -1 - \lambda_0 - \lambda_1(x - \mu)^2$$

$$q(x) = e^{-\lambda_0 - 1 - \lambda_1(x - \mu)^2}$$

- This is a "shape" of searched pdf
- Now we need to find the Lagrangian multipliers by:

$$\frac{d\mathcal{L}(q(x))}{d\lambda_0} = 0$$

$$\frac{d\mathcal{L}(q(x))}{d\lambda_1} = 0$$

Alternatively we can substitute the found shape to each constraint (a) and (b)

Continuous R.Vs.: max entropy (app 1)

Proof: max entropy for bounded variance

- Satisfy the first constraint:

$$\int q(x)dx = 1 = \int e^{-\lambda_0 - 1 - \lambda_1(x-\mu)^2} dx$$

$$1 = \int e^{-\lambda_0 - 1 - \lambda_1 z^2} dz \mid \text{Change of variables } z = x - \mu$$

$$1 = \int e^{-\lambda_0 - 1} e^{-\lambda_1 z^2} dz$$

$$e^{\lambda_0 + 1} = \int e^{-\lambda_1 z^2} dz$$

$$e^{\lambda_0 + 1} = \sqrt{\frac{\pi}{\lambda_1}}$$

Continuous R.Vs.: max entropy (app 1)

- Satisfy the second constraint:

$$\int q(x)(x - \mu)^2 dx = \sigma^2 = \int e^{-\lambda_0 - 1 - \lambda_1(x - \mu)^2} (x - \mu)^2 dx$$

$$\sigma^2 = \int e^{-\lambda_0 - 1 - \lambda_1(x - \mu)^2} (x - \mu)^2 dx$$

$$\sigma^2 = \int e^{-\lambda_0 - 1 - \lambda_1 z^2} z^2 dz \mid \text{Change of variables } z = x - \mu$$

$$\sigma^2 e^{\lambda_0 + 1} = \frac{1}{2} \sqrt{\frac{\pi}{\lambda_1^3}}$$

$$\sigma^2 e^{\lambda_0 + 1} = \frac{1}{2\lambda_1} \sqrt{\frac{\pi}{\lambda_1}}$$

$$2\lambda_1 \sigma^2 e^{\lambda_0 + 1} = \sqrt{\frac{\pi}{\lambda_1}}$$

Continuous R.Vs.: max entropy (app 1)

- Putting together

$$\sqrt{\frac{\pi}{\lambda_1}} = e^{\lambda_0+1} = 2\lambda_1\sigma^2 e^{\lambda_0+1}$$

$$e^{\lambda_0-1} = 2\lambda_1\sigma^2 e^{\lambda_0-1}$$

$$1 = 2\lambda_1\sigma^2$$

$$\lambda_1 = \frac{1}{2\sigma^2}$$

Continuous R.Vs.: max entropy (app 1)

- Putting together

$$\sqrt{\frac{\pi}{\lambda_1}} = e^{\lambda_0+1} = 2\lambda_1\sigma^2e^{\lambda_0+1}$$

$$e^{\lambda_0-1} = 2\lambda_1\sigma^2e^{\lambda_0-1}$$

$$1 = 2\lambda_1\sigma^2$$

$$\lambda_1 = \frac{1}{2\sigma^2}$$

$$\sqrt{\frac{\pi}{\lambda_1}} = e^{\lambda_0+1}$$

$$\sqrt{2\sigma^2\pi} = e^{\lambda_0+1}$$

$$\ln \sqrt{2\sigma^2\pi} = \lambda_0 + 1$$

$$\lambda_0 = \ln \sqrt{2\sigma^2\pi} - 1$$

Continuous R.Vs.: max entropy (app 1)

- Substituting $\lambda_1 = \frac{1}{2\sigma^2}$
 $\lambda_0 = \ln \sqrt{2\sigma^2\pi} - 1$
 to $q(x) = e^{-\lambda_0 - 1 - \lambda_1(x-\mu)^2}$ results

$$q(x) = e^{-\ln \sqrt{2\sigma^2\pi} + 1 - 1 - \frac{1}{2\sigma^2}(x-\mu)^2}$$

$$q(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \Rightarrow \mathcal{N}(\mu, \sigma_X^2)$$

- Thus, the distribution possessing the maximum entropy among all distributions with the bounded variance is Gaussian.

Continuous R.Vs.: max entropy for fixed var (app 2)

Proof: max entropy for bounded variance

- Gaussian has maximum differential entropy **among all distributions with the same variance**
- Let $q(x)$ be any density satisfying $\int x^2 q(x) dx = \sigma^2$. Let $p(x) = \mathcal{N}(0, \sigma_X^2)$. Then

$$h(q(x)) \leq h(p(x))$$

$$\begin{aligned} 0 &\leq D_{\text{KL}}(q(x) \| p(x)) = \int q(x) \log \frac{q(x)}{p(x)} dx \quad \text{We will see all KLD properties later} \\ &= -h(q(x)) - \int q(x) \log p(x) dx \\ &= -h(q(x)) - \int p(x) \log p(x) dx \\ &= -h(q(x)) + h(p(x)) \end{aligned}$$

$$\begin{aligned} \log p(x) &= -\log(\sqrt{2\sigma^2}) - x^2/2\sigma^2 \\ \int q(x)x^2 dx &= \int p(x)x^2 dx = \sigma^2 \end{aligned}$$

$$h(q(x)) \leq h(p(x))$$

Continuous R.Vs.: max entropy for fixed var (app 2)

- Demo

$$0 \leq D_{\text{KL}}(q(x) \| p(x)) = \int q(x) \log \frac{q(x)}{p(x)} dx$$

$$-D_{\text{KL}}(q(x) \| p(x)) = \int q(x) \log \frac{p(x)}{q(x)} dx$$

$$\leq \log \int q(x) \frac{p(x)}{q(x)} dx \text{ (by Jensen's inequality)}$$

$$= \log \int p(x) dx \leq \log = 0$$

Multivariate Gaussian: Differential entropy

Definition (Differential entropy of Gaussian r.v.): *differential entropy of Gaussian random vector* $\mathbf{X} \sim \mathcal{N}(\bar{\mathbf{x}}, \mathbf{K}_{\mathbf{xx}})$ is:

$$h(\mathbf{X}) = -E_{f_{\mathbf{X}}} [\log_2 f_{\mathbf{X}}(\mathbf{x})] = \frac{1}{2} \ln \left((2\pi e)^N |\det \mathbf{K}_{\mathbf{xx}}| \right)$$

- Proof

$$f(\mathbf{x}) = f(x_1, x_2, \dots, x_N) = \frac{1}{\sqrt{(2\pi)^N |\det \mathbf{K}_{\mathbf{xx}}|}} \exp \left[-\frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}})^T \mathbf{K}_{\mathbf{xx}}^{-1} (\mathbf{x} - \bar{\mathbf{x}}) \right]$$

$$h(\mathbf{X}) = - \int_{\mathbb{R}^N} f(\mathbf{x}) \ln f(\mathbf{x}) d\mathbf{x} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f(\mathbf{x}) \ln f(\mathbf{x}) d\mathbf{x} =$$

Multivariate Gaussian: Differential entropy

$$\begin{aligned} h(\mathbf{X}) &= - \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f(\mathbf{x}) \ln f(\mathbf{x}) d\mathbf{x} \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} f(\mathbf{x}) \left[\frac{1}{2} \ln \left((2\pi)^N |\det \mathbf{K}_{\mathbf{xx}}| \right) + \frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}})^T \mathbf{K}_{\mathbf{xx}}^{-1} (\mathbf{x} - \bar{\mathbf{x}}) \right] d\mathbf{x} \\ &= \frac{1}{2} \ln \left((2\pi)^N |\det \mathbf{K}_{\mathbf{xx}}| \right) + \frac{1}{2} \underbrace{E_{f(\mathbf{x})} \left[(\mathbf{x} - \bar{\mathbf{x}})^T \mathbf{K}_{\mathbf{xx}}^{-1} (\mathbf{x} - \bar{\mathbf{x}}) \right]}_N \ln e \\ &= \frac{1}{2} \ln \left((2\pi)^N |\det \mathbf{K}_{\mathbf{xx}}| \right) + \frac{1}{2} N \ln e \\ &= \frac{1}{2} \ln \left((2\pi e)^N |\det \mathbf{K}_{\mathbf{xx}}| \right) \end{aligned}$$

- We will prove that $E_{f(\mathbf{x})} \left[(\mathbf{x} - \bar{\mathbf{x}})^T \mathbf{K}_{\mathbf{xx}}^{-1} (\mathbf{x} - \bar{\mathbf{x}}) \right] = N$

Multivariate Gaussian: Differential entropy

- Note $\underbrace{(\mathbf{x} - \bar{\mathbf{x}})^T}_{1 \times N} \underbrace{\mathbf{K}_{\mathbf{xx}}^{-1}}_{N \times N} \underbrace{(\mathbf{x} - \bar{\mathbf{x}})}_{N \times 1} \Rightarrow 1 \times 1$
- Since the trace of any 1×1 matrix is $Tr[\mathbf{A}] = \mathbf{A}$, one obtains

$$(\mathbf{x} - \bar{\mathbf{x}})^T \mathbf{K}_{\mathbf{xx}}^{-1} (\mathbf{x} - \bar{\mathbf{x}}) = Tr[(\mathbf{x} - \bar{\mathbf{x}})^T \mathbf{K}_{\mathbf{xx}}^{-1} (\mathbf{x} - \bar{\mathbf{x}})]$$
- Thus

$$\begin{aligned} E_{f(\mathbf{x})} \left[(\mathbf{x} - \bar{\mathbf{x}})^T \mathbf{K}_{\mathbf{xx}}^{-1} (\mathbf{x} - \bar{\mathbf{x}}) \right] &= \int_{-\infty}^{+\infty} f(\mathbf{x}) (\mathbf{x} - \bar{\mathbf{x}})^T \mathbf{K}_{\mathbf{xx}}^{-1} (\mathbf{x} - \bar{\mathbf{x}}) d\mathbf{x} \\ &= \int_{-\infty}^{+\infty} f(\mathbf{x}) \ln e^{(\mathbf{x}-\bar{\mathbf{x}})^T \mathbf{K}_{\mathbf{xx}}^{-1} (\mathbf{x}-\bar{\mathbf{x}})} d\mathbf{x} \\ &= \int_{-\infty}^{+\infty} f(\mathbf{x}) \ln e^{Tr[(\mathbf{x}-\bar{\mathbf{x}})^T \mathbf{K}_{\mathbf{xx}}^{-1} (\mathbf{x}-\bar{\mathbf{x}})]} d\mathbf{x} \end{aligned}$$

Multivariate Gaussian: Differential entropy

- Note $\text{Tr}[\mathbf{ABC}] = \text{Tr}[\mathbf{CAB}] = \text{Tr}[\mathbf{BCA}]$

- Thus

$$\int_{-\infty}^{+\infty} f(\mathbf{x}) \ln e^{\text{Tr}[(\mathbf{x}-\bar{\mathbf{x}})^T \mathbf{K}_{\mathbf{xx}}^{-1} (\mathbf{x}-\bar{\mathbf{x}})]} d\mathbf{x} = \int_{-\infty}^{+\infty} f(\mathbf{x}) \ln e^{\text{Tr}[\mathbf{K}_{\mathbf{xx}}^{-1} (\mathbf{x}-\bar{\mathbf{x}})(\mathbf{x}-\bar{\mathbf{x}})^T]} d\mathbf{x}$$

- Because $f(\mathbf{x})$ is a scalar and the natural logarithm and exponential may cancel, the properties of the trace function allow us to push $f(\mathbf{x})$ and the integral inside the trace

$$\begin{aligned} \int_{-\infty}^{+\infty} f(\mathbf{x}) \ln e^{\text{Tr}[\mathbf{K}_{\mathbf{xx}}^{-1} (\mathbf{x}-\bar{\mathbf{x}})(\mathbf{x}-\bar{\mathbf{x}})^T]} d\mathbf{x} &= \ln e^{\underbrace{\text{Tr}[\mathbf{K}_{\mathbf{xx}}^{-1} \int_{-\infty}^{+\infty} f(\mathbf{x})(\mathbf{x}-\bar{\mathbf{x}})(\mathbf{x}-\bar{\mathbf{x}})^T d\mathbf{x}]}_{\mathbf{K}_{\mathbf{xx}}}} \\ &= \ln e^{\text{Tr}[\mathbf{K}_{\mathbf{xx}}^{-1} \mathbf{K}_{\mathbf{xx}}]} \\ &= \ln e^{\text{Tr}[\mathbf{I}_N]} = \ln e^N = N \end{aligned}$$

Multivariate Gaussian: Differential entropy

Definition (Differential entropy of Gaussian r.v.): *differential entropy of Gaussian random vector* $\mathbf{X} \sim \mathcal{N}(\bar{\mathbf{x}}, \mathbf{K}_{\mathbf{xx}})$ is:

$$h(\mathbf{X}) = \frac{1}{2} \ln \left((2\pi e)^N \det |\mathbf{K}_{\mathbf{xx}}| \right) = \frac{1}{2} \ln \left((2\pi e)^N \prod_{i=1}^N \sigma_i \right) = \frac{1}{2} \sum_{i=1}^N \ln (2\pi e \sigma_i)$$

$$\mathbf{K}_{\mathbf{xx}} = \mathbf{U}\Sigma\mathbf{U}^T \Rightarrow \left| \det \mathbf{K}_{\mathbf{xx}} \right| = \left| \det \mathbf{U}\Sigma\mathbf{U}^T \right| = \underbrace{\left| \det \mathbf{U}\mathbf{U}^T \right|}_{\mathbf{I}} \left| \det \Sigma \right| = \left| \det \Sigma \right| = \prod_{i=1}^N \sigma_i$$

$\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n)$ or $\Sigma_{ii} = \sigma_i, \forall i$ - eigenvalues of covariance matrix (not variances)

$$\mathbf{K}_{x_i x_j} = \text{Cov}[X_i, X_j] = \mathbb{E}[(X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])] = \sigma_{X_{ij}}$$

$$\mathbf{K}_{x_i x_i} = \text{Cov}[X_i, X_i] = \mathbb{E}[(X_i - \mathbb{E}[X_i])(X_i - \mathbb{E}[X_i])] = \sigma_{X_i}^2 \text{ - variances}$$

Multivariate Gaussian: Differential entropy

Application (Hadamard's inequality) *differential entropy of Gaussian random vector $\mathbf{X} \sim \mathcal{N}(\bar{\mathbf{x}}, \mathbf{K}_{\mathbf{xx}})$* can be bounded as:

$$h(\mathbf{X}) = \frac{1}{2} \ln((2\pi e)^N |\det \mathbf{K}_{\mathbf{xx}}|) \leq \frac{1}{2} \ln \left((2\pi e)^N \prod_{i=1}^N \sigma_{X_i}^2 \right) = \frac{1}{2} \sum_{i=1}^N \ln(2\pi e \sigma_{X_i}^2)$$

- Hadamard's inequality $|\mathbf{K}_{\mathbf{xx}}| \leq \prod_{i=1}^N \mathbf{K}_{x_i x_i} = \prod_{i=1}^N \sigma_{X_i}^2$

$$\mathbf{K}_{\mathbf{xx}} = \begin{bmatrix} \sigma_{X_1}^2 & \rho \sigma_{X_1} \sigma_{X_2} \\ \rho \sigma_{X_1} \sigma_{X_2} & \sigma_{X_2}^2 \end{bmatrix}$$

$$|\det \mathbf{K}_{\mathbf{xx}}| = \sigma_{X_1}^2 \sigma_{X_2}^2 - \underbrace{\rho^2}_{0 \leq \rho^2 \leq 1} \sigma_{X_1}^2 \sigma_{X_2}^2 \leq \sigma_{X_1}^2 \sigma_{X_2}^2 = \prod_{i=1}^N \sigma_{X_i}^2$$

Multivariate Gaussian: Differential entropy

Application (Hadamard's inequality) *differential entropy of Gaussian random vector $\mathbf{X} \sim \mathcal{N}(\bar{\mathbf{x}}, \mathbf{K}_{\mathbf{xx}})$ can be bounded as:*

$$h(\mathbf{X}) = \frac{1}{2} \ln((2\pi e)^N |\det \mathbf{K}_{\mathbf{xx}}|) \leq \frac{1}{2} \ln \left((2\pi e)^N \prod_{i=1}^N \sigma_{X_i}^2 \right) = \frac{1}{2} \sum_{i=1}^N \ln(2\pi e \sigma_{X_i}^2)$$

- Alternative “IT proof”

$$\underbrace{\frac{1}{2} \ln((2\pi e)^N |\det \mathbf{K}_{\mathbf{xx}}|)}_{\text{general case}} = h(\mathbf{X}) \leq \underbrace{\sum_{i=1}^N h(X_i)}_{\text{independent variables}} = \frac{1}{2} \sum_{i=1}^N \ln(2\pi e \sigma_{X_i}^2)$$

Continuous R.Vs.: Differential entropy

- Translation does not change entropy

$$h(X + a) = h(X)$$

- Impact of scaling on differential entropy

- Scalar case

$$h(aX) = h(X) + \log|a|$$

$$\begin{aligned} X &\sim \mathcal{N}(0, \sigma_X^2) \\ aX &\sim \mathcal{N}(0, a^2 \sigma_X^2) \end{aligned}$$

$$h(aX) = \frac{1}{2} \log_2 (2\pi e a^2 \sigma_X^2) = \frac{1}{2} \log_2 (2\pi e \sigma_X^2) + \log_2 |a|$$

- Vector case

$$h(\mathbf{AX}) = h(\mathbf{X}) + \log|\det(\mathbf{A})|$$

Continuous R.Vs.: Differential entropy

- Impact of scaling on differential entropy
 - Formal proof

$$h(aX) = h(X) + \log|a|$$

$$Y = aX \quad f_Y(y) = \frac{1}{|a|} f_X(y/a)$$

$$h(aX) = - \int f_Y(y) \log f_Y(y) dy = - \int \frac{1}{|a|} f_X(y/a) \log \left(\frac{1}{|a|} f_X(y/a) \right) dy$$

$$= - \int \frac{1}{|a|} f_X(y/a) \underbrace{\log \left(\frac{1}{|a|} \right)}_{\text{const}} dy - \underbrace{\int \frac{1}{|a|} f_X(y/a) \log f_X(y/a) dy}_{h(X)}$$

$$= h(X) + \log|a|$$

Summary: entropy of discrete vs continuous RVs



- Entropy of discrete RVs

- Negativity

$$H(X) \geq 0$$

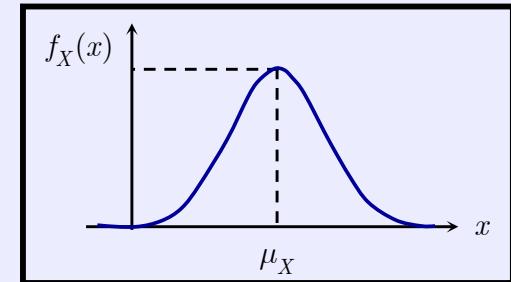
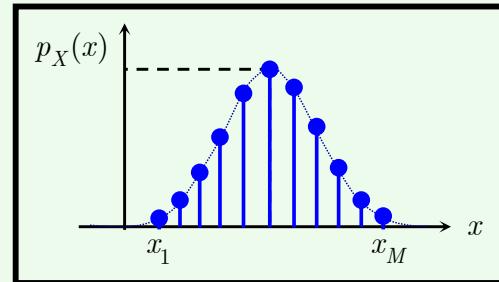
- Scaling

$$H(aX) = H(X)$$

- Entropy of continuous RVs

$$h(X) \text{ any}$$

$$h(aX) = h(X) + \log|a|$$



Continuous R.Vs.: max entropy for fixed cov matrix

- Let the random vector $\mathbf{X} \in \mathbb{R}^N$ has zero mean and covariance $\mathbf{K}_{\mathbf{xx}} = E[\mathbf{XX}^T]$, i.e., $\mathbf{K}_{x_i x_j} = E[X_i X_j]$, $1 \leq i, j \leq N$.
- Then any random vector with the above covariance matrix has the entropy smaller or equal to the Gaussian one $h(\mathbf{X}) \leq \frac{1}{2} \ln((2\pi e)^N |\mathbf{K}_{\mathbf{xx}}|)$

- Let $q(\mathbf{x})$ be any density satisfying $\int \mathbf{x}_i \mathbf{x}_j q(\mathbf{x}) d\mathbf{x} = E_q[X_i X_j] = \mathbf{K}_{x_i x_j}$
- Let $p(\mathbf{x})$ be the density of a $\mathcal{N}(\mathbf{0}, \mathbf{K}_{\mathbf{xx}})$, i.e., it has $\int \mathbf{x}_i \mathbf{x}_j p(\mathbf{x}) d\mathbf{x} = \mathbf{K}_{x_i x_j}$
 - Note that $\log p(\mathbf{x}) = -\log \left(\sqrt{(2\pi)^N |\det \mathbf{K}_{\mathbf{xx}}|} \right) - \left[-\frac{1}{2} \mathbf{x}^T \mathbf{K}_{\mathbf{xx}}^{-1} \mathbf{x} \right]$

Continuous R.Vs.: max entropy for fixed cov matrix

$$0 \leq D_{\text{KL}}(q(\mathbf{x}) \| p(\mathbf{x})) = \int q(\mathbf{x}) \log \frac{q(\mathbf{x})}{p(\mathbf{x})} d\mathbf{x}$$

$$= -h(q(\mathbf{x})) - \int q(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x}$$

$$= -h(q(\mathbf{x})) - \int p(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x}$$

$$= -h(q(\mathbf{x})) + h(p(\mathbf{x}))$$

$$h(q(\mathbf{x})) \leq h(p(\mathbf{x}))$$

$$\int q(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x} = \int p(\mathbf{x}) \log p(\mathbf{x}) d\mathbf{x}$$



$$\log p(\mathbf{x}) = -\log \left(\sqrt{(2\pi)^N |\det \mathbf{K}_{\mathbf{xx}}|} \right) - \left[-\frac{1}{2} \mathbf{x}^T \mathbf{K}_{\mathbf{xx}}^{-1} \mathbf{x} \right]$$

$$\int q(\mathbf{x}) \mathbf{x}^T \mathbf{K}_{\mathbf{xx}}^{-1} \mathbf{x} d\mathbf{x} = N$$

$$\int p(\mathbf{x}) \mathbf{x}^T \mathbf{K}_{\mathbf{xx}}^{-1} \mathbf{x} d\mathbf{x} = N$$

See “Proof of expectation of quadratic form”

Continuous R.V.: relative entropy

□
Definition (relative entropy): *Relative entropy* or *Kullback-Leibler divergence* between pdfs $p(x)$ and $q(x)$:

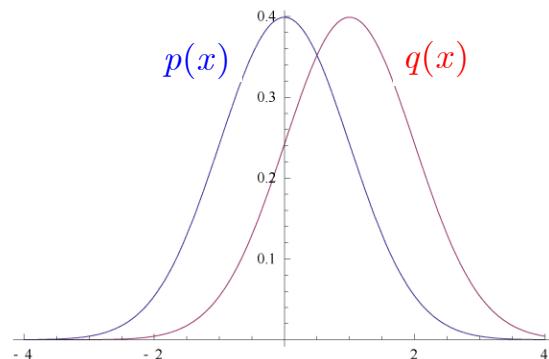
$$D_{\text{KL}}(p \parallel q) = \int_{-\infty}^{\infty} p(x) \log \left(\frac{p(x)}{q(x)} \right) dx = E_{p(x)} \left[\log_2 \frac{p(x)}{q(x)} \right]$$

- Equivalent notations

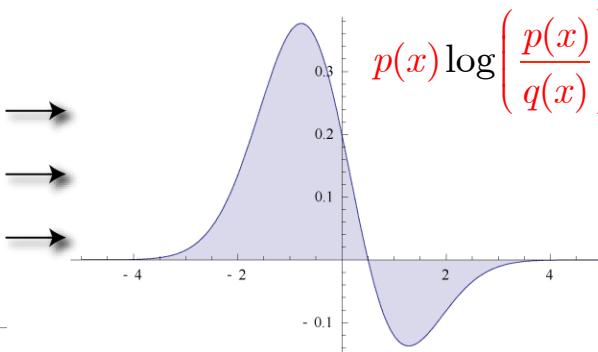
$$D_{\text{KL}}(f \parallel q) = \int_{\mathcal{X}} f(x) \log_2 \frac{f(x)}{q(x)} dx$$

Continuous R.V.: relative entropy

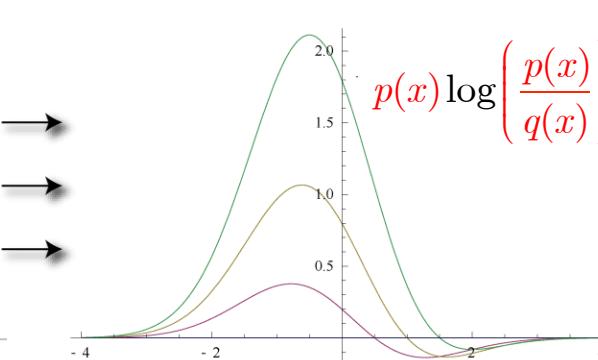
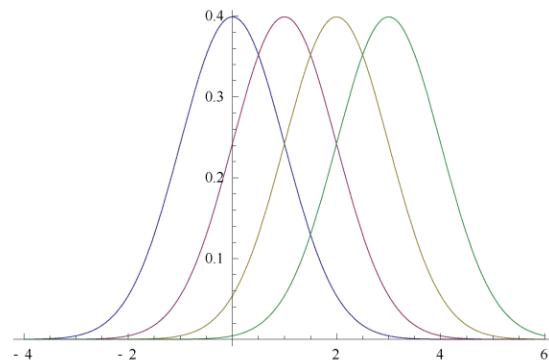
$$D_{\text{KL}}(p \parallel q) = \int_{-\infty}^{\infty} p(x) \log \left(\frac{p(x)}{q(x)} \right) dx = \underbrace{E_{p(x)} [\log_2 p(x)]}_{-h(p(x))} - \underbrace{E_{p(x)} [\log_2 q(x)]}_{h(p(x); q(x))}$$



Original Gaussian PDF's



KL Area to be Integrated $D_{\text{KL}}(p \parallel q)$



Continuous R.Vs.: KLD between two Gaussians

□
Definition (relative entropy between two Gaussian): *Relative entropy* between two univariate Gaussian pdfs $p(x) = \mathcal{N}(\mu_p, \sigma_p^2)$ and $q(x) = \mathcal{N}(\mu_q, \sigma_q^2)$

$$\begin{aligned} D_{\text{KL}}(p \parallel q) &= D_{\text{KL}}(\mathcal{N}(\mu_p, \sigma_p^2) \parallel \mathcal{N}(\mu_q, \sigma_q^2)) \\ &= \log \sigma_q - \log \sigma_p - \frac{1}{2} \left[1 - \left(\frac{\sigma_p^2 + (\mu_p - \mu_q)^2}{\sigma_q^2} \right) \right] \end{aligned}$$

- Proof

$$D_{\text{KL}}(p \parallel q) = \int_{-\infty}^{\infty} p(x) \log \left(\frac{p(x)}{q(x)} \right) dx = E_{p(x)} \left[\log_2 \frac{p(x)}{q(x)} \right]$$

Continuous R.Vs.: KLD between two Gaussians

$$\log \frac{p(x)}{q(x)} = \log \left(\frac{\frac{1}{\sqrt{2\pi\sigma_p^2}} \exp \left[-\frac{1}{2\sigma_p^2} (x - \mu_p)^2 \right]}{\frac{1}{\sqrt{2\pi\sigma_q^2}} \exp \left[-\frac{1}{2\sigma_q^2} (x - \mu_q)^2 \right]} \right)$$

$$= \log \left(\frac{\sqrt{2\pi\sigma_q^2}}{\sqrt{2\pi\sigma_p^2}} + \frac{1}{2\sigma_q^2} (x - \mu_q)^2 - \frac{1}{2\sigma_p^2} (x - \mu_p)^2 \right)$$

$$D_{\text{KL}}(p \parallel q) = E_{p(x)} \left[\log_2 \frac{p(x)}{q(x)} \right] = \frac{1}{2} \log \left(\frac{\sigma_q^2}{\sigma_p^2} \right) + \frac{1}{2\sigma_q^2} E_{p(x)} \left[(X - \mu_q)^2 \right]$$

$$- \frac{1}{2\sigma_p^2} \underbrace{E_{p(x)} \left[(X - \mu_p)^2 \right]}_{\sigma_p^2}$$

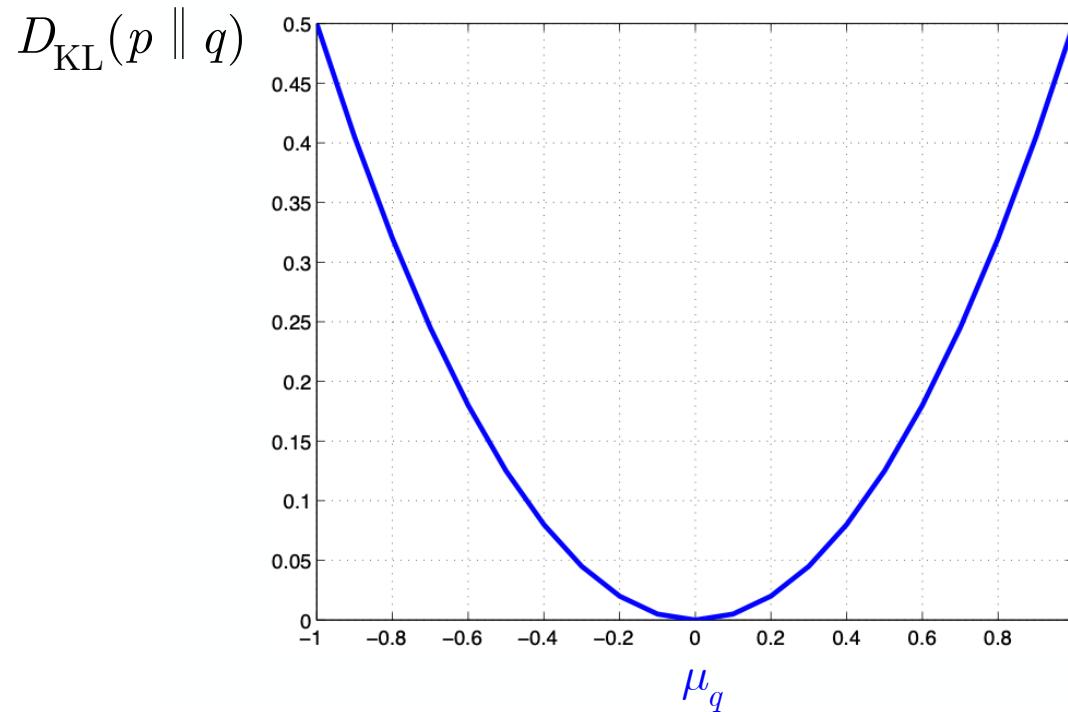
Continuous R.Vs.: KLD between two Gaussians

$$\begin{aligned}
 E_{p(x)}[(X - \mu_q)^2] &= E_{p(x)}[X^2 - 2X\mu_q + \mu_q^2] \\
 &= \underbrace{E_{p(x)}[X^2]}_{\sigma_p^2 + \mu_p^2} - 2\mu_q \underbrace{E_{p(x)}[X]}_{\mu_p} + \mu_q^2 \\
 &= \sigma_p^2 + \mu_p^2 - 2\mu_q\mu_p + \mu_q^2 \\
 &= \sigma_p^2 - (\mu_p - \mu_q)^2
 \end{aligned}$$

$$\begin{aligned}
 D_{\text{KL}}(p \parallel q) &= \frac{1}{2} \log \left(\frac{\sigma_q^2}{\sigma_p^2} \right) + \frac{1}{2\sigma_q^2} E_{p(x)}[(X - \mu_q)^2] - \frac{1}{2\sigma_p^2} E_{p(x)}[(X - \mu_p)^2] \\
 &= \log \sigma_q - \log \sigma_p + \frac{1}{2\sigma_q^2} [\sigma_p^2 - (\mu_p - \mu_q)^2] - \frac{1}{2\sigma_p^2} \sigma_p^2 \\
 &= \log \sigma_q - \log \sigma_p - \frac{1}{2} \left[1 - \left(\frac{\sigma_p^2 + (\mu_p - \mu_q)^2}{\sigma_q^2} \right) \right]
 \end{aligned}$$

Continuous R.Vs.: KLD between two Gaussians

$$p(x) = \mathcal{N}(0,1) \quad q(x) = \mathcal{N}(\mu_q, 1)$$



Continuous R.V.: multivariate relative entropy

Definition (relative entropy): *Relative entropy* or *Kullback-Leibler divergence*

between pdfs $p(\mathbf{x})$ and $q(\mathbf{x})$:

$$D_{\text{KL}}(p \parallel q) = \int_{\mathbb{R}^N} p(\mathbf{x}) \log \left(\frac{p(\mathbf{x})}{q(\mathbf{x})} \right) d\mathbf{x} = E_{p(\mathbf{x})} \left[\log_2 \frac{p(\mathbf{x})}{q(\mathbf{x})} \right]$$

Continuous R.V.: KLD for multivariate Gaussians

Definition (relative entropy): *Kullback-Leibler divergence* between two multivariate Gaussian pdfs $p(\mathbf{x}) = \mathcal{N}(\bar{\mathbf{x}}_p, \mathbf{K}_p)$ and $q(\mathbf{x}) = \mathcal{N}(\bar{\mathbf{x}}_q, \mathbf{K}_q)$:

$$D_{\text{KL}}(p(\mathbf{x}) \parallel q(\mathbf{x})) = \frac{1}{2} \left(\ln \left(\frac{\det \mathbf{K}_q}{\det \mathbf{K}_p} \right) + \text{tr} \left(\mathbf{K}_q^{-1} \mathbf{K}_p \right) + (\bar{\mathbf{x}}_q - \bar{\mathbf{x}}_p)^T \mathbf{K}_q^{-1} (\bar{\mathbf{x}}_q - \bar{\mathbf{x}}_p) - N \right)$$

- Proof: we will consider two approaches to prove it
 - Approach A: direct proof
 - Approach B: proof via entropy-cross entropy decomposition

Continuous R.V.: KLD for multivariate Gaussians

- Recall $p(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^N |\det \mathbf{K}_p|}} \exp\left[-\frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}}_p)^T \mathbf{K}_p^{-1} (\mathbf{x} - \bar{\mathbf{x}}_p)\right]$
- $\ln p(\mathbf{x}) = -\frac{1}{2} \ln\left((2\pi)^N |\det \mathbf{K}_p|\right) - \frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}}_p)^T \mathbf{K}_p^{-1} (\mathbf{x} - \bar{\mathbf{x}}_p)$
- Substitute $p(\mathbf{x}) = \mathcal{N}(\bar{\mathbf{x}}_p, \mathbf{K}_p)$ and $q(\mathbf{x}) = \mathcal{N}(\bar{\mathbf{x}}_q, \mathbf{K}_q)$

$$\begin{aligned}
 D_{\text{KL}}(p(\mathbf{x}) \parallel q(\mathbf{x})) &= E_{p(\mathbf{x})} [\ln p(\mathbf{x}) - \ln q(\mathbf{x})] \\
 &= \frac{1}{2} E_{p(\mathbf{x})} \left[-\ln(|\det \mathbf{K}_p|) - (\mathbf{x} - \bar{\mathbf{x}}_p)^T \mathbf{K}_p^{-1} (\mathbf{x} - \bar{\mathbf{x}}_p) + \ln(|\det \mathbf{K}_q|) + (\mathbf{x} - \bar{\mathbf{x}}_q)^T \mathbf{K}_q^{-1} (\mathbf{x} - \bar{\mathbf{x}}_q) \right] \\
 &= \frac{1}{2} \ln \left(\frac{|\det \mathbf{K}_q|}{|\det \mathbf{K}_p|} \right) + \frac{1}{2} E_{p(\mathbf{x})} \left[-(\mathbf{x} - \bar{\mathbf{x}}_p)^T \mathbf{K}_p^{-1} (\mathbf{x} - \bar{\mathbf{x}}_p) + (\mathbf{x} - \bar{\mathbf{x}}_q)^T \mathbf{K}_q^{-1} (\mathbf{x} - \bar{\mathbf{x}}_q) \right] \\
 &= \frac{1}{2} \ln \left(\frac{|\det \mathbf{K}_q|}{|\det \mathbf{K}_p|} \right) + \frac{1}{2} E_{p(\mathbf{x})} \left[-Tr[\mathbf{K}_p^{-1} (\mathbf{x} - \bar{\mathbf{x}}_p) (\mathbf{x} - \bar{\mathbf{x}}_p)^T] + Tr[\mathbf{K}_q^{-1} (\mathbf{x} - \bar{\mathbf{x}}_q) (\mathbf{x} - \bar{\mathbf{x}}_q)^T] \right]
 \end{aligned}$$

Continuous R.V.: KLD for multivariate Gaussians

$$\dots = \frac{1}{2} \ln \left(\frac{\left| \det \mathbf{K}_q \right|}{\left| \det \mathbf{K}_p \right|} \right) + \frac{1}{2} \mathbb{E}_{p(\mathbf{x})} \left[-Tr \left[\mathbf{K}_p^{-1} (\mathbf{x} - \bar{\mathbf{x}}_p) (\mathbf{x} - \bar{\mathbf{x}}_p)^T \right] + Tr \left[\mathbf{K}_q^{-1} (\mathbf{x} - \bar{\mathbf{x}}_q) (\mathbf{x} - \bar{\mathbf{x}}_q)^T \right] \right]$$

$$\begin{aligned} \mathbb{E}_{p(\mathbf{x})} \left[-Tr \left[\mathbf{K}_p^{-1} (\mathbf{x} - \bar{\mathbf{x}}_p) (\mathbf{x} - \bar{\mathbf{x}}_p)^T \right] \right] &= -Tr \left[\mathbf{K}_p^{-1} \mathbb{E}_{p(\mathbf{x})} \left[(\mathbf{x} - \bar{\mathbf{x}}_p) (\mathbf{x} - \bar{\mathbf{x}}_p)^T \right] \right] \\ &= -Tr \left[\mathbf{K}_p^{-1} \mathbf{K}_p \right] = -Tr \left[\mathbf{I}_N \right] = -N \end{aligned}$$

$$\begin{aligned} \mathbb{E}_{p(\mathbf{x})} \left[Tr \left[\mathbf{K}_q^{-1} (\mathbf{x} - \bar{\mathbf{x}}_q) (\mathbf{x} - \bar{\mathbf{x}}_q)^T \right] \right] &= \mathbb{E}_{p(\mathbf{x})} \left[Tr \left[\mathbf{K}_q^{-1} \left[\mathbf{x}\mathbf{x}^T - 2\mathbf{x}\bar{\mathbf{x}}_q^T + \bar{\mathbf{x}}_q\bar{\mathbf{x}}_q^T \right] \right] \right] \\ &= Tr \left[\mathbf{K}_q^{-1} \left[\underbrace{\mathbb{E}_{p(\mathbf{x})} \left[\mathbf{x}\mathbf{x}^T \right]}_{\mathbf{K}_p + \bar{\mathbf{x}}_p\bar{\mathbf{x}}_p^T} - 2\underbrace{\mathbb{E}_{p(\mathbf{x})} \left[\mathbf{x} \right] \bar{\mathbf{x}}_q^T}_{\bar{\mathbf{x}}_p} + \bar{\mathbf{x}}_q\bar{\mathbf{x}}_q^T \right] \right] \\ &= Tr \left[\mathbf{K}_q^{-1} \left[\mathbf{K}_p + \bar{\mathbf{x}}_p\bar{\mathbf{x}}_p^T - 2\bar{\mathbf{x}}_p\bar{\mathbf{x}}_q^T + \bar{\mathbf{x}}_q\bar{\mathbf{x}}_q^T \right] \right] \\ &= Tr \left[\mathbf{K}_q^{-1} \left[\mathbf{K}_p + \bar{\mathbf{x}}_p\bar{\mathbf{x}}_p^T - 2\bar{\mathbf{x}}_p\bar{\mathbf{x}}_q^T + \bar{\mathbf{x}}_q\bar{\mathbf{x}}_q^T \right] \right] \\ &= Tr \left[\mathbf{K}_q^{-1} \mathbf{K}_p \right] + Tr \left[\mathbf{K}_q^{-1} \left[\bar{\mathbf{x}}_p\bar{\mathbf{x}}_p^T - 2\bar{\mathbf{x}}_p\bar{\mathbf{x}}_q^T + \bar{\mathbf{x}}_q\bar{\mathbf{x}}_q^T \right] \right] \\ &= Tr \left[\mathbf{K}_q^{-1} \mathbf{K}_p \right] + Tr \left[\underbrace{\bar{\mathbf{x}}_p \mathbf{K}_q^{-1} \bar{\mathbf{x}}_p^T - 2\bar{\mathbf{x}}_p \mathbf{K}_q^{-1} \bar{\mathbf{x}}_q^T + \bar{\mathbf{x}}_q \mathbf{K}_q^{-1} \bar{\mathbf{x}}_q^T}_{(\bar{\mathbf{x}}_q - \bar{\mathbf{x}}_p)^T \mathbf{K}_q^{-1} (\bar{\mathbf{x}}_q - \bar{\mathbf{x}}_p)} \right] \\ &= Tr \left[\mathbf{K}_q^{-1} \mathbf{K}_p \right] + (\bar{\mathbf{x}}_q - \bar{\mathbf{x}}_p)^T \mathbf{K}_q^{-1} (\bar{\mathbf{x}}_q - \bar{\mathbf{x}}_p) \end{aligned}$$

Continuous R.V.: KLD for multivariate Gaussians

- Recall $D_{\text{KL}}(p \parallel q) = h(p, q) - h(p)$

$$h(p) = -E_{p(\mathbf{x})}[\ln p(\mathbf{x})] = \frac{1}{2} \ln \left((2\pi e)^N \det |\mathbf{K}_p| \right) = \frac{1}{2} \ln \left((2\pi)^N \det |\mathbf{K}_p| \right) + \underbrace{\frac{1}{2} \ln(e^N)}_N$$

$$h(p, q) = -E_{p(\mathbf{x})}[\ln q(\mathbf{x})]$$

$$\ln q(\mathbf{x}) = -\frac{1}{2} \ln \left((2\pi)^N |\det \mathbf{K}_q| \right) - \frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}}_q)^T \mathbf{K}_q^{-1} (\mathbf{x} - \bar{\mathbf{x}}_q)$$

$$h(p, q) = -E_{p(\mathbf{x})} \left[-\frac{1}{2} \ln \left((2\pi)^N |\det \mathbf{K}_q| \right) - \frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}}_q)^T \mathbf{K}_q^{-1} (\mathbf{x} - \bar{\mathbf{x}}_q) \right]$$

$$= \frac{1}{2} \ln \left((2\pi)^N |\det \mathbf{K}_q| \right) + \underbrace{E_{p(\mathbf{x})} \left[\frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}}_q)^T \mathbf{K}_q^{-1} (\mathbf{x} - \bar{\mathbf{x}}_q) \right]}_{Tr[\mathbf{K}_q^{-1} \mathbf{K}_p] + (\bar{\mathbf{x}}_q - \bar{\mathbf{x}}_p)^T \mathbf{K}_q^{-1} (\bar{\mathbf{x}}_q - \bar{\mathbf{x}}_p)}$$

Continuous R.V.: multivariate relative entropy

Relative entropy between $p(\mathbf{x}) = \mathcal{N}(\bar{\mathbf{x}}_p, \Sigma_p)$ and $q(\mathbf{x}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$

$$D_{\text{KL}}\left(\mathcal{N}\left(\left(\mu_1, \dots, \mu_N\right)^T, \text{diag}\left(\sigma_1^2, \dots, \sigma_N^2\right)\right) \| \mathcal{N}(\mathbf{0}, \mathbf{I})\right) = \frac{1}{2} \sum_{i=1}^N \left(\sigma_i^2 + \mu_i^2 - \ln(\sigma_i^2) - 1 \right)$$

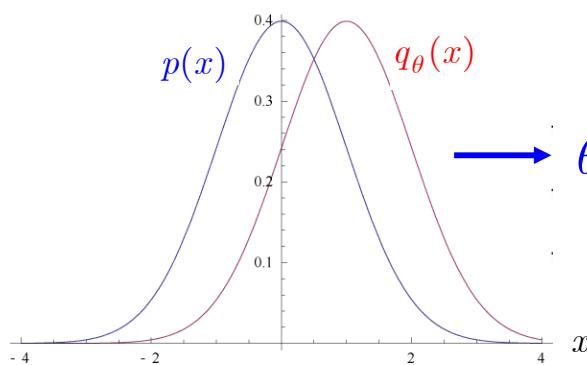
$$\bar{\mathbf{x}}_p = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_N \end{bmatrix} \quad \Sigma_p = \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_N^2 \end{bmatrix} \quad \bar{\mathbf{x}}_q = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix} \quad \Sigma_q = \begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix}$$

KLD in action

- We consider properties of KLD on several ML applications
 - Approximation of true pdf by approximate pdf
 - Binary hypothesis testing

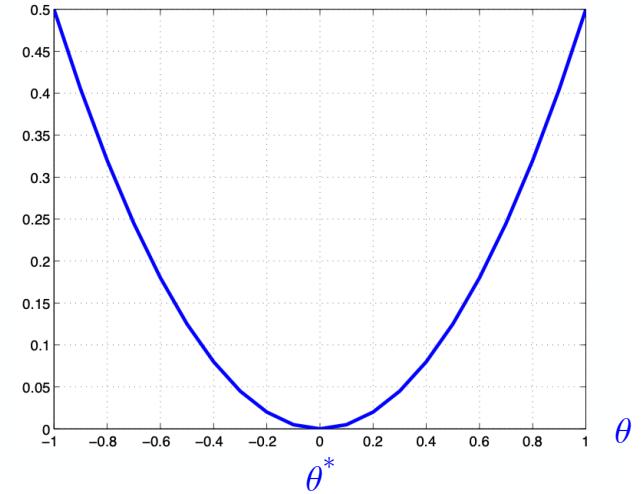
Continuous R.V.: KLD for approximation

- We consider properties of KLD on one major ML application
- Approximation of a given true distribution $p(x)$ by some approximative distribution $q_\theta(x)$ controlled by a set of parameters θ



$$D_{\text{KL}}(p \parallel q_\theta)$$

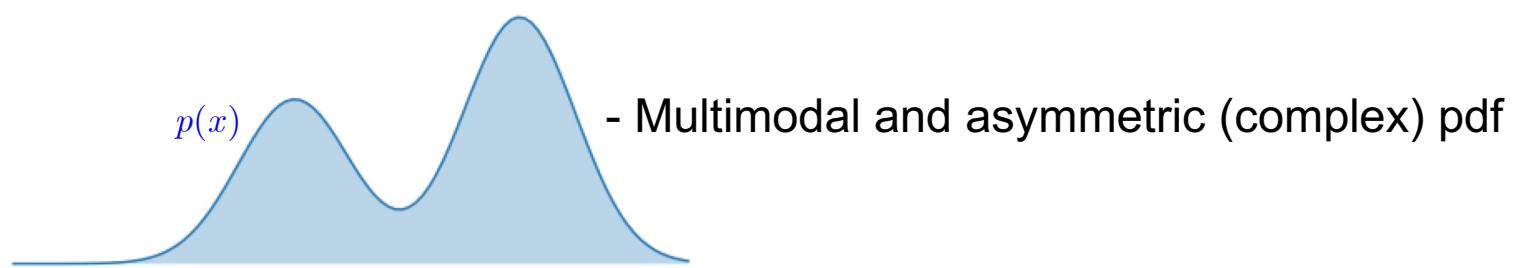
$$\theta^* = \arg \min_{\theta} D_{\text{KL}}(p \parallel q_\theta)$$



- KLD is used as a “measure” of closeness of approximation
- Our problem: should we use $D_{\text{KL}}(p \parallel q_\theta)$ or $D_{\text{KL}}(q_\theta \parallel p)$?
 - Recall $D_{\text{KL}}(p \parallel q_\theta) \neq D_{\text{KL}}(q_\theta \parallel p)$

Continuous R.V.: Forward and Reverse KLD

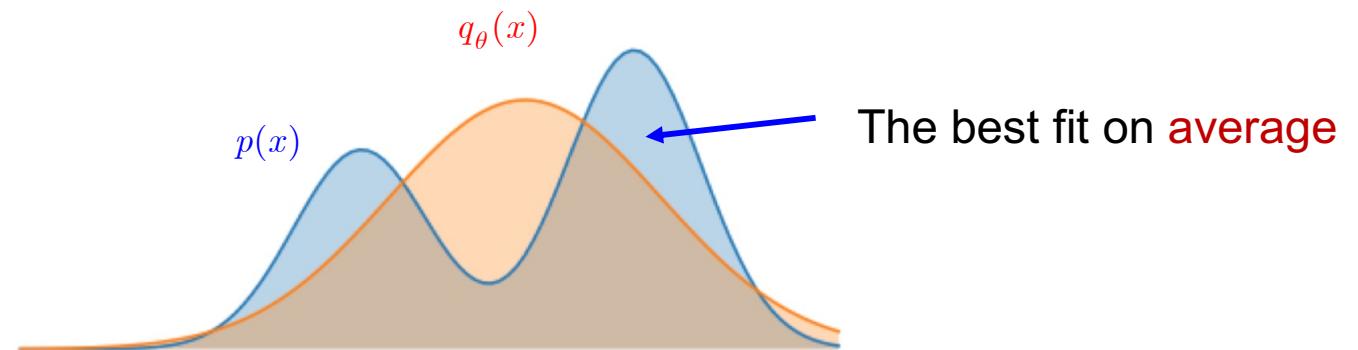
- We have two options:
 - **Forward KLD:** $\arg \min_{\theta} D_{\text{KL}}(p \parallel q_{\theta})$
 - **Reverse KLD:** $\arg \min_{\theta} D_{\text{KL}}(q_{\theta} \parallel p)$
- Example: let $q(x) = \mathcal{N}(\mu, \sigma^2); \theta = (\mu, \sigma^2)$ - unimodal (simple) pdf



Continuous R.V.: Forward and Reverse KLD

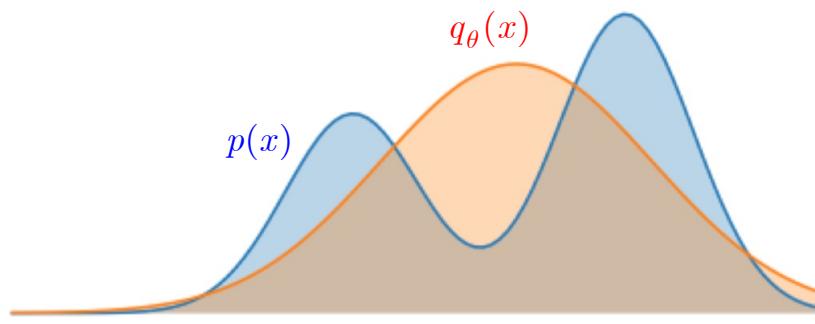
- **Forward KL: Mean-Seeking Behavior**

$$\begin{aligned}\arg \min_{\theta} D_{\text{KL}}(p \parallel q_{\theta}) &= \arg \min_{\theta} \underbrace{E_{p(x)}[\log p(x)]}_{-H(p(x)) \text{ const wrt } \theta} - \underbrace{E_{p(x)}[\log q_{\theta}(x)]}_{+H(p(x);q_{\theta}(x))} \\ &\equiv \arg \min_{\theta} \underbrace{-E_{p(x)}[\log q_{\theta}(x)]}_{H(p(x);q_{\theta}(x))} \\ &= \arg \max_{\theta} \underbrace{E_{p(x)}[\log q_{\theta}(x)]}_{H(p(x);q_{\theta}(x))} \\ &= \arg \min_{\theta} H(p(x);q_{\theta}(x))\end{aligned}$$



Continuous R.V.: Forward and Reverse KLD

- **Forward KL: Mean-Seeking Behavior**

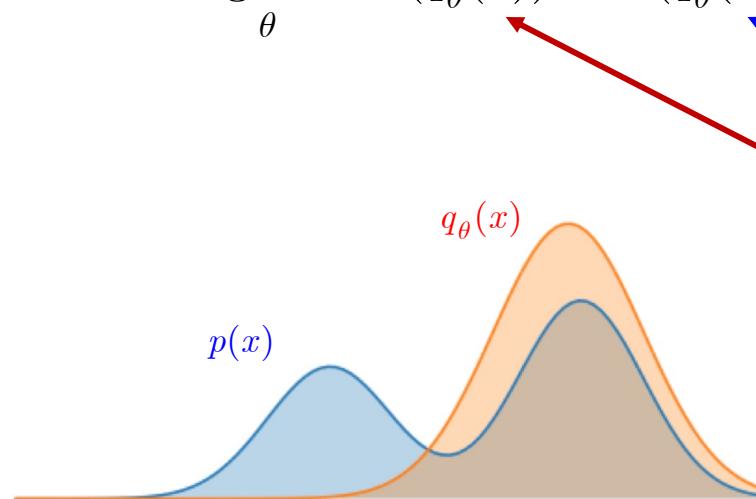


- The approximate distribution centers itself between the two modes, so that it can have high coverage of both
- The maximum (mode) of approximate distribution does not coincide with the mode of true distribution

Continuous R.V.: Forward and Reverse KLD

- Reverse KL: Mode-Seeking Behavior

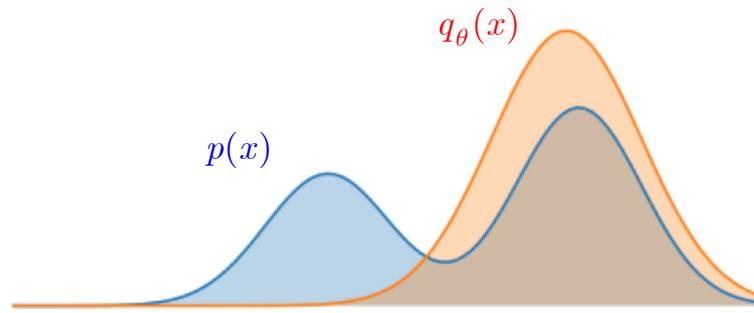
$$\begin{aligned}\arg \min_{\theta} D_{\text{KL}}(q_{\theta} \| p) &= \arg \min_{\theta} \underbrace{E_{q_{\theta}(x)}[\log q_{\theta}(x)]}_{-H(q_{\theta}(x))} \underbrace{-E_{q_{\theta}(x)}[\log p(x)]}_{+H(q_{\theta}(x); p(x))} \\ &\equiv \arg \min_{\theta} -H(q_{\theta}(x)) + H(q_{\theta}(x); p(x)) \\ &= \arg \max_{\theta} H(q_{\theta}(x)) - H(q_{\theta}(x); p(x))\end{aligned}$$



The best fit with the regularization that the searched distribution should have max entropy at the same time

Continuous R.V.: Forward and Reverse KLD

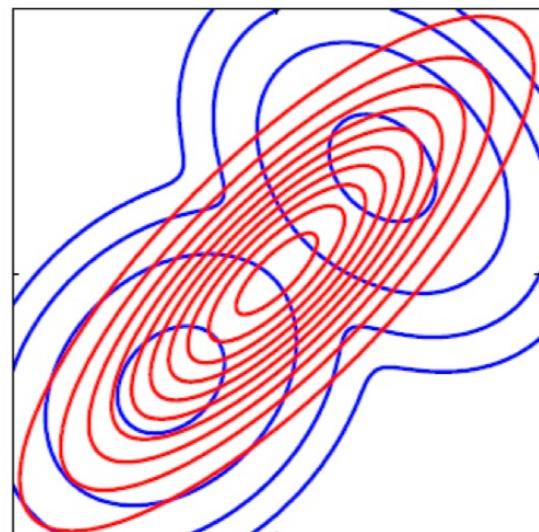
- Reverse KL: **Mode-Seeking** Behavior



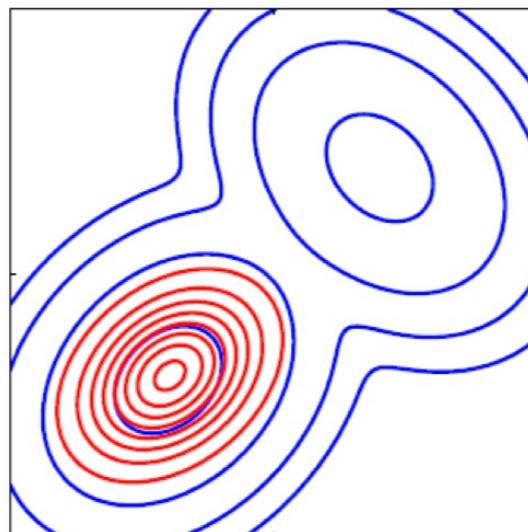
- The approximate distribution centers itself at the mode
- At the same time, the cross-entropy term prevents the approximate distribution from collapsing to a very narrow mode
- Thus, optimizing this objective is to find a mode of true distribution with high probability and wide support and mimic true distribution

Continuous R.V.: Forward and Reverse KLD

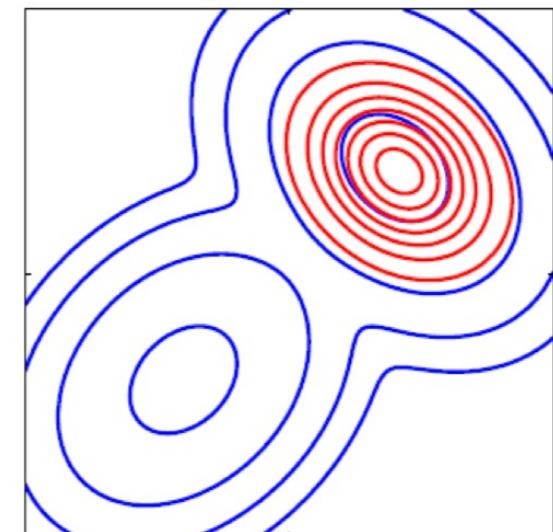
- Forward and reverse KLD in 2D



(a)



(b)



(c)

$$\arg \min_{\theta} D_{\text{KL}}(p \parallel q_{\theta})$$

$$\arg \min_{\theta} D_{\text{KL}}(q_{\theta} \parallel p)$$

Figure from *Bishop, Pattern Recognition and Machine Learning, 2006*

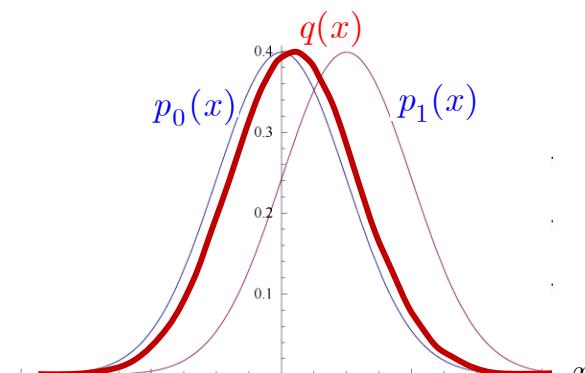
KLD in action

- We consider properties of KLD on several ML applications
 - Approximation of true pdf by approximate pdf
 - Binary hypothesis testing (HT)

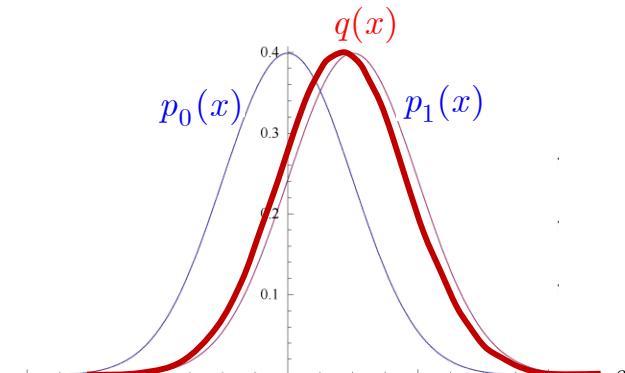
Continuous R.V.: KLD for binary HT

- We consider the use of KLD for binary sequential hypothesis testing

- i.i.d.*
- Given a sequence of observation $X_1, X_2, \dots, X_N \sim q(x)$
 - Given two models $H_0 : p_0(x)$ and $H_1 : p_1(x)$, we should decide from which distribution the above sequence is generated from or “belongs” to
 - Intuition: what is the best fit $D_{\text{KL}}(q \parallel p_0)$ or $D_{\text{KL}}(q \parallel p_1)$?



$$D_{\text{KL}}(q \parallel p_0) < D_{\text{KL}}(q \parallel p_1) \Rightarrow H_0$$



$$D_{\text{KL}}(q \parallel p_0) > D_{\text{KL}}(q \parallel p_1) \Rightarrow H_1$$

Continuous R.V.: KLD for binary HT

- If $q(x)$ is known, we know how to do it *i.i.d.*
- However, we only know samples $X_1, X_2, \dots, X_N \sim q(x)$ generated from $q(x)$
- Consider likelihood ratio test (LRT) for N independent observations

$$\mathcal{L}(X_1, X_2, \dots, X_N) = \prod_{i=1}^N \frac{p_1(x_i)}{p_0(x_i)} \stackrel{H_1}{\gtrless} \gamma \stackrel{H_0}{\gtrless}$$

- Normalized (divided by N) log-likelihood ratio test

$$\ell(X_1, X_2, \dots, X_N) = \frac{1}{N} \log \mathcal{L}(X_1, X_2, \dots, X_N) = \frac{1}{N} \sum_{i=1}^N \log \frac{p_1(x_i)}{p_0(x_i)} \stackrel{H_1}{\gtrless} \ln \gamma \stackrel{H_0}{\gtrless}$$

Continuous R.V.: KLD for binary HT

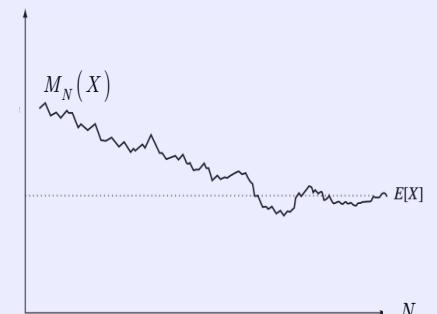
- **Remark:** $\ell(X_1, X_2, \dots, X_N)$ is a random variable representing a sum
- From the strong law of large numbers

- Let X_1, X_2, \dots, X_N be a sequence of i.i.d. random variables with finite mean $E[X_i] = \mu_X$ and variance $Var[X_i] = \sigma_X^2$ (some bounded)

- In this case

$$\lim_{N \rightarrow \infty} \Pr\left[\lim_{N \rightarrow \infty} M_N(X) = \mu_X \right] = 1$$

- $M_K(X)$ is the sample mean computed using K variables



$$\frac{1}{N} \sum_{i=1}^N \log \frac{p_1(x_i)}{p_0(x_i)} \xrightarrow{\text{as } N \rightarrow \infty} E_{q(x)} \left[\log \frac{p_1(x)}{p_0(x)} \right]$$

Continuous R.V.: KLD for binary HT

$$\begin{aligned} E_{q(x)} \left[\log \frac{p_1(x)}{p_0(x)} \right] &= E_{q(x)} \left[\log \frac{p_1(x)}{p_0(x)} \frac{q(x)}{q(x)} \right] \\ &= E_{q(x)} \left[\log \frac{q(x)}{p_0(x)} \right] - E_{q(x)} \left[\log \frac{q(x)}{p_1(x)} \right] \\ &= D_{\text{KL}}(q \parallel p_0) - D_{\text{KL}}(q \parallel p_1) \end{aligned}$$

$$\frac{1}{N} \sum_{i=1}^N \log \frac{p_1(x_i)}{p_0(x_i)} \stackrel{H_1}{\gtrless} \ln \gamma \quad \Rightarrow \quad D_{\text{KL}}(q \parallel p_0) - D_{\text{KL}}(q \parallel p_1) \stackrel{H_1}{\gtrless} \ln \gamma$$

- **Remark:** for equal probabilities of hypothesis and minimum of average

probability of errors $\ln \gamma = 0$: $\stackrel{H_1}{\gtrless}$

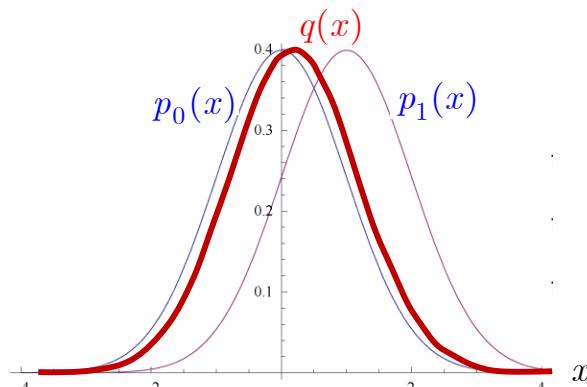
$$D_{\text{KL}}(q \parallel p_0) \stackrel{H_0}{\gtrless} D_{\text{KL}}(q \parallel p_1)$$

Continuous R.V.: KLD for binary HT

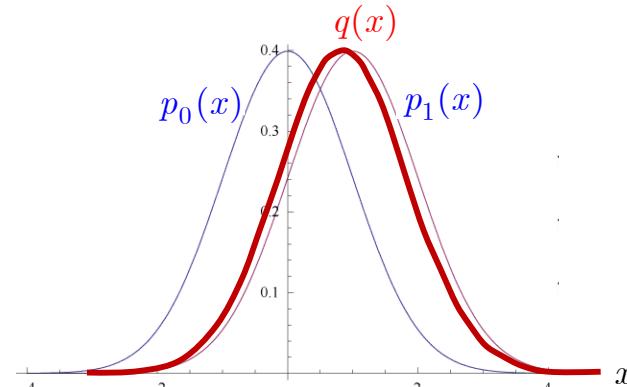
- Conclusion

$$\frac{1}{N} \sum_{i=1}^N \log \frac{p_1(x_i)}{p_0(x_i)} \stackrel{H_1}{\gtrless} \ln \gamma \Rightarrow D_{\text{KL}}(q \parallel p_0) - D_{\text{KL}}(q \parallel p_1) \stackrel{H_1}{\gtrless} \ln \gamma$$

$$D_{\text{KL}}(q \parallel p_0) \stackrel{H_1}{\gtrless} D_{\text{KL}}(q \parallel p_1)$$



$$D_{\text{KL}}(q \parallel p_0) < D_{\text{KL}}(q \parallel p_1) \Rightarrow H_0$$



$$D_{\text{KL}}(q \parallel p_0) > D_{\text{KL}}(q \parallel p_1) \Rightarrow H_1$$

Continuous R.V.: conditional relative entropy

Definition (conditional relative entropy): *Conditional relative entropy* between pdfs $p(y|x)$ and $q(y|x)$:

$$\begin{aligned} D_{\text{KL}}(p(y|x) \| q(y|x)) &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} p(x)p(y|x) \log \left(\frac{p(y|x)}{q(y|x)} \right) dx dy \\ &= \underbrace{E_{p(x)} E_{p(y|x)}}_{E_{p(x,y)}} \left[\log_2 \frac{p(y|x)}{q(y|x)} \right] \\ &= E_{p(x,y)} \left[\log_2 \frac{p(y|x)}{q(y|x)} \right] \end{aligned}$$

Continuous R.V.: chain rule for relative entropy

- Chain rule

$$D_{\text{KL}}(p(x, y) \parallel q(x, y)) = D_{\text{KL}}(p(x) \parallel q(x)) + D_{\text{KL}}(p(y|x) \parallel q(y|x))$$

$$\begin{aligned} &= \int_{x \in \mathcal{X}} \int_{y \in \mathcal{Y}} p(y, x) \log_2 \frac{p(y|x)p(x)}{q(y|x)q(x)} dx dy = \\ &= \int_{x \in \mathcal{X}} \int_{y \in \mathcal{Y}} p(y, x) \left(\log_2 \frac{p(y|x)}{q(y|x)} + \log_2 \frac{p(x)}{q(x)} \right) dx dy = \\ &= \int_{x \in \mathcal{X}} \int_{y \in \mathcal{Y}} p(y, x) \log_2 \frac{p(y|x)}{q(y|x)} dx dy + \int_{x \in \mathcal{X}} \int_{y \in \mathcal{Y}} p(y, x) \log_2 \frac{p(x)}{q(x)} dx dy = \\ &= D_{\text{KL}}(p(x) \parallel q(x)) + D_{\text{KL}}(p(y|x) \parallel q(y|x)) \end{aligned}$$

Continuous R.Vs.: Mutual information

□
Definition (muual information): *Mutual information* for continuous random variable has the same properties as for the discrete random variablea

$$\begin{aligned} I(X;Y) &= \int_x \int_y f_{X,Y}(x,y) \log_2 \frac{f_{X,Y}(x,y)}{f_X(x)f_Y(y)} dx dy \\ &= D_{\text{KL}}(f_{X,Y}(x,y) \| f_X(x)f_Y(y)) \end{aligned}$$

$$I(X;Y) \geq 0$$

$$I(X;Y) = h(X) - h(X|Y)$$

$$I(X;Y) = h(Y) - h(Y|X)$$

$$I(X;Y) = h(X) + h(Y) - h(X,Y)$$

Continuous R.Vs.: Mutual information

$$\begin{aligned} I(X;Y) &= D_{\text{KL}}(f_{X,Y}(x,y) \parallel f_X(x)f_Y(y)) \\ &= E_{f_Y(y)} \left[D_{\text{KL}}(f_{X|Y=y}(x|y) \parallel f_X(x)) \right] = E_{f_X(x)} \left[D_{\text{KL}}(f_{Y|X=x}(y|x) \parallel f_Y(y)) \right] \end{aligned}$$

- Proof

$$\begin{aligned} I(X;Y) &= \underbrace{\sum_{y \in \mathcal{Y}} f_Y(y) \sum_{x \in \mathcal{X}} f_{X|Y=y}(x|y) \log \frac{f_{X|Y=y}(x|y)}{f_X(x)}}_{\sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} f_{X,Y}(x,y)} \\ &= \sum_{y \in \mathcal{Y}} f_Y(y) D_{\text{KL}}(f_{X|Y=y}(x|y) \parallel f_X(x)) \\ &= E_{f_Y(y)} \left[D_{\text{KL}}(f_{X|Y=y}(x|y) \parallel f_X(x)) \right] \end{aligned}$$

Remark: here we highlight the fact of conditioning in the KLD. However, many ML works simplify it to the notations on the previous pages.

Continuous R.Vs.: Mutual information

- Conditional mutual information

$$I(X;Y|Z) = H(X|Z) - H(X|Y,Z)$$

- Proof

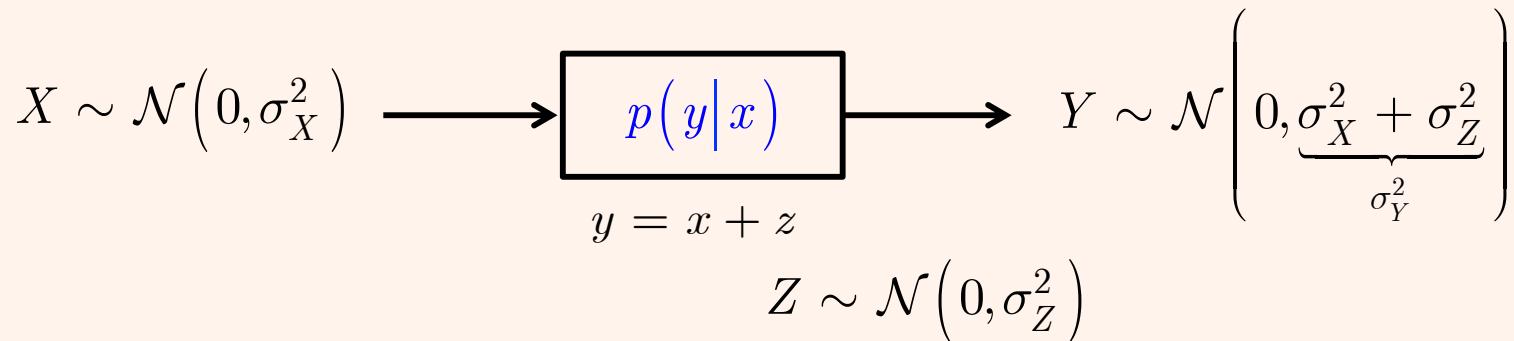
$$\begin{aligned} I(X;Y|Z) &= E_{f(x,y,z)} \left[\log \frac{f_{X,Y|Z}(x,y|z)}{f_{X|Z}(x|z)f_{Y|Z}(y|z)} \right] = E_{f(x,y,z)} \left[\log \frac{f_{X|Y,Z}(x|y,z)}{f_{X|Z}(x|z)} \right] \\ f_{X,Y|Z}(x,y|z) &= f_{Y|Z}(y|z)f_{X|Y,Z}(x|y,z) \\ &= \underbrace{E_{f(x,y,z)} \left[\log f_{X|Y,Z}(x|y,z) \right]}_{-h(X|Y,Z)} - \underbrace{E_{f(x,y,z)} \left[\log f_{X|Z}(x|z) \right]}_{h(X|Z)} \end{aligned}$$

- Chain rule

$$I(X_1, X_2, \dots, X_N; Y) = \sum_{i=1}^N I(X_i; Y | X_{i-1}, \dots, X_1)$$

Continuous R.Vs.: MI for univariate Gaussians

- Consider a transformation of Gaussian random variable $X \sim \mathcal{N}(0, \sigma_X^2)$ by adding an independent random variable $Z \sim \mathcal{N}(0, \sigma_Z^2)$, i.e., $Y = X + Z$



- Compute mutual information via entropy and show the equivalence of all forms

$$I(X;Y) = h(Y) - h(Y|X) = \frac{1}{2} \log \left(1 + \frac{\sigma_X^2}{\sigma_Z^2} \right)$$

$$I(X;Y) = h(X) - h(X|Y) = \frac{1}{2} \log \left(1 + \frac{\sigma_X^2}{\sigma_Z^2} \right) \quad I(X;Y) = D_{\text{KL}}(p(x,y) \| p(x)p(y))$$

$$I(X;Y) = h(X) + h(Y) - h(X,Y) = \frac{1}{2} \log \left(1 + \frac{\sigma_X^2}{\sigma_Z^2} \right)$$

Continuous R.Vs.: MI for univariate Gaussians

$$I(X;Y) = h(Y) - h(Y|X) = \frac{1}{2} \log \left(1 + \frac{\sigma_X^2}{\sigma_Z^2} \right)$$

- Proof

- Random variable $Y \sim \mathcal{N} \left(0, \underbrace{\sigma_X^2 + \sigma_Z^2}_{\sigma_Y^2} \right)$

$$h(Y) = \frac{1}{2} \log(2\pi e \sigma_Y^2) = \frac{1}{2} \log(2\pi e (\sigma_X^2 + \sigma_Z^2))$$

- Since $Y = X + Z$

$$h(Y|X) = h(X+Z|X) = h(Z) = \frac{1}{2} \log(2\pi e \sigma_Z^2)$$

- Finally

$$I(X;Y) = h(Y) - h(Y|X) = \frac{1}{2} \log \left(\frac{2\pi e (\sigma_X^2 + \sigma_Z^2)}{2\pi e \sigma_Z^2} \right) = \frac{1}{2} \log \left(1 + \frac{\sigma_X^2}{\sigma_Z^2} \right)$$

Continuous R.Vs.: MI for univariate Gaussians

$$I(X;Y) = h(X) - h(X|Y) = \frac{1}{2} \log \left(1 + \frac{\sigma_X^2}{\sigma_Z^2} \right)$$

- Proof

- Random variable $X \sim \mathcal{N}(0, \sigma_X^2)$

$$h(X) = \frac{1}{2} \log(2\pi e \sigma_X^2) = \frac{1}{2} \log(2\pi e \sigma_X^2)$$

- Since $Y = X + Z$

$$h(X|Y) = h(X|X+Z) = \frac{1}{2} \log(2\pi e \sigma_{X|Y}^2)$$

where $\sigma_{X|Y}^2 = \frac{\sigma_X^2 \sigma_Z^2}{\sigma_X^2 + \sigma_Z^2}$

- Finally

$$I(X;Y) = h(X) - h(X|Y) = \frac{1}{2} \log \left(\frac{2\pi e(\sigma_X^2 + \sigma_Z^2)}{2\pi e \sigma_Z^2} \right) = \frac{1}{2} \log \left(1 + \frac{\sigma_X^2}{\sigma_Z^2} \right)$$

Continuous R.Vs.: MI for univariate Gaussians

Proof: MAP estimation of signal in noise

- Estimation of X given $Y = X + Z$ based on **maximum a posteriori (MAP)** principle

$$\hat{x} = \arg \max_x p(y|x)p(x)$$

$$p(y|x) = \mathcal{N}(x, \sigma_Z^2) = \frac{1}{\sqrt{2\pi\sigma_Z^2}} e^{-\frac{(y-x)^2}{2\sigma_Z^2}} \quad p(x) = \mathcal{N}(0, \sigma_X^2) = \frac{1}{\sqrt{2\pi\sigma_X^2}} e^{-\frac{(x-0)^2}{2\sigma_X^2}}$$

$$\hat{x} = \arg \max_x \log[p(y|x)p(x)] \Rightarrow \frac{d}{dx} [\log p(y|x) + \log p(x)] = 0$$

$$\frac{d}{dx} \left[-\frac{(y-x)^2}{2\sigma_Z^2} - \frac{(x)^2}{2\sigma_X^2} \right] = 0 \Rightarrow \frac{2}{2\sigma_Z^2}(y-x) - \frac{2x}{2\sigma_X^2} = 0$$

$$\frac{y}{\sigma_Z^2} - \left(\frac{1}{\sigma_Z^2} + \frac{1}{\sigma_X^2} \right)x = 0 \Rightarrow \hat{x} = \frac{\sigma_Z^2 \sigma_X^2}{\sigma_Z^2 + \sigma_X^2} \frac{1}{\sigma_Z^2} y = \frac{\sigma_X^2}{\sigma_Z^2 + \sigma_X^2} y$$

Continuous R.Vs.: MI for univariate Gaussians

- Variance of MAP estimation

$$\sigma_{X|Y}^2 = E_{p(x)} \left[|X - \hat{X}|^2 \right] = \frac{\sigma_X^2 \sigma_Z^2}{\sigma_X^2 + \sigma_Z^2}$$

- Proof

$$\begin{aligned}
 \sigma_{X|Y}^2 &= E_{p(x,z)} \left[|X - \hat{X}|^2 \right] = E_{p(x,z)} \left[\left| X - \frac{\sigma_X^2}{\sigma_X^2 + \sigma_Z^2} Y \right|^2 \right] \\
 &= E_{p(x,z)} \left[\left| X - \frac{\sigma_X^2}{\sigma_X^2 + \sigma_Z^2} (X + Z) \right|^2 \right] = E_{p(x,z)} \left[\left| \frac{\sigma_Z^2}{\sigma_X^2 + \sigma_Z^2} X - \frac{\sigma_X^2}{\sigma_X^2 + \sigma_Z^2} Z \right|^2 \right] \\
 &= \left(\frac{\sigma_Z^2}{\sigma_X^2 + \sigma_Z^2} \right)^2 \sigma_X^2 + \left(\frac{\sigma_X^2}{\sigma_X^2 + \sigma_Z^2} \right)^2 \sigma_Z^2 = \frac{\sigma_X^2 \sigma_Z^2 (\sigma_Z^2 + \sigma_X^2)}{(\sigma_X^2 + \sigma_Z^2)^2} = \frac{\sigma_X^2 \sigma_Z^2}{\sigma_X^2 + \sigma_Z^2}
 \end{aligned}$$

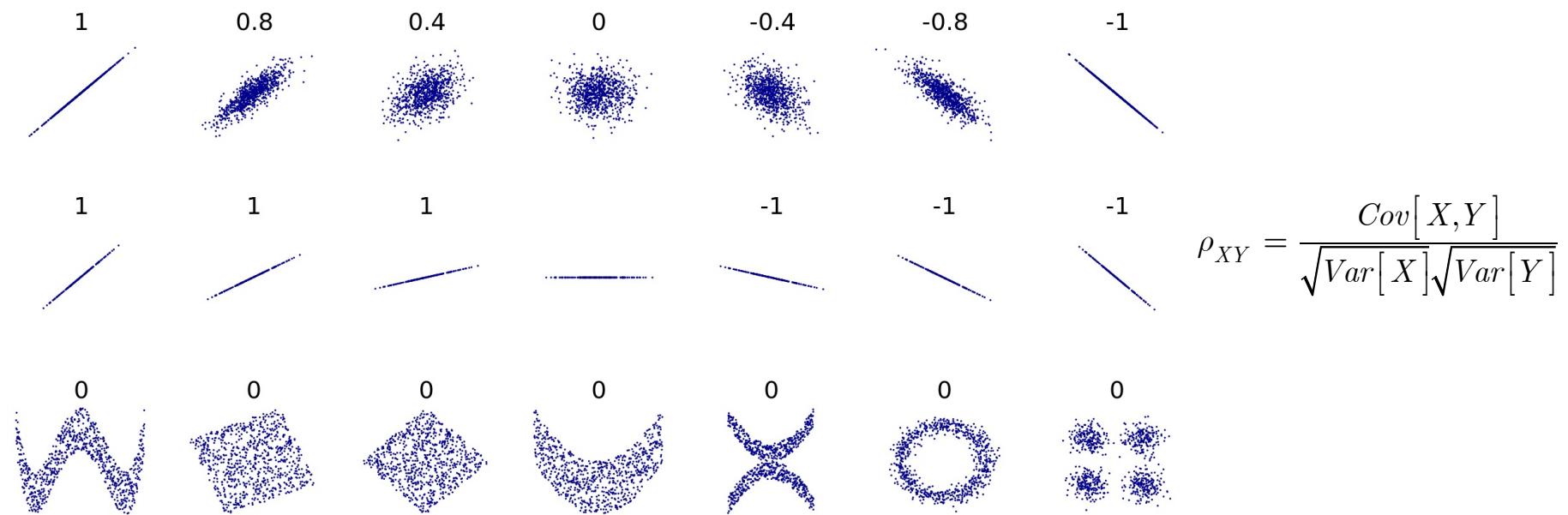
Continuous R.Vs.: MI for univariate Gaussians

$$I(X;Y) = h(X) + h(Y) - h(X,Y) = \frac{1}{2} \log \left(1 + \frac{\sigma_X^2}{\sigma_Z^2} \right) = \frac{1}{2} \log \left(\frac{1}{1 - \rho_{XY}^2} \right)$$

- Remark
 - This exercise is to show not only the equivalence of three forms of mutual information decomposition via the entropy but also to demonstrate:
 - How to compute and use the joint entropy
 - How to link it to the correlation coefficient
 - Note: mutual information vs correlation coefficient (or MSE)
 - Mutual information contains information about **all** dependencies – linear and nonlinear
 - Correlation coefficient measures **only** linear dependencies
 - Note: in the bivariate Gaussian case between X and Y they coincide

Continuous R.Vs.: MI for univariate Gaussians

- Correlation coefficient for different data manifolds (pdfs)



- Note that the correlation reflects the strength and direction of a linear relationship (top row), but not the slope of that relationship (middle), nor many aspects of nonlinear relationships (bottom).
- N.B.: the figure in the center has a slope of 0 but in that case the correlation coefficient is undefined because the variance of Y is zero.

Continuous R.Vs.: MI for univariate Gaussians

$$I(X;Y) = h(X) + h(Y) - h(X,Y) = \frac{1}{2} \log\left(1 + \frac{\sigma_X^2}{\sigma_Z^2}\right) = \frac{1}{2} \log\left(\frac{1}{1 - \rho_{XY}^2}\right)$$

- Proof

- Random variable $X \sim \mathcal{N}(\mu_X, \sigma_X^2)$ and $Y \sim \mathcal{N}(\mu_Y, \sigma_Y^2)$

$$h(X) = \frac{1}{2} \log(2\pi e \sigma_X^2)$$

$$h(Y) = \frac{1}{2} \log(2\pi e \sigma_Y^2)$$

- To find $h(X, Y)$, we will define a joint pdf:

$(X, Y) \sim p(x, y) = \mathcal{N}(\underline{\mu}, \mathbf{K}_{xy})$ which is Gaussian since $Y = X + Z$

- where

$$\underline{\mu} = \begin{bmatrix} \mu_X \\ \mu_Y \end{bmatrix} \quad \mathbf{K}_{xy} = \begin{bmatrix} \sigma_X^2 & \rho_{XY} \sigma_X \sigma_Y \\ \rho_{XY} \sigma_X \sigma_Y & \sigma_Y^2 \end{bmatrix}$$

Continuous R.Vs.: MI for univariate Gaussians

- With the correlation coefficient

$$\begin{aligned}
 \rho_{XY} &= \frac{\text{Cov}[X, Y]}{\sqrt{\text{Var}[X]}\sqrt{\text{Var}[Y]}} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y} = \frac{E[XY] - \mu_X \mu_Y}{\sigma_X \sigma_Y} \\
 &= \frac{E[X(X + Z)] - \mu_X(\mu_X + 0)}{\sigma_X \sigma_Y} = \frac{E[X^2] - \mu_X^2}{\sigma_X \sigma_Y} \\
 &= \frac{\sigma_X^2}{\sigma_X \sigma_Y} = \frac{\sigma_X}{\sigma_Y} = \pm \sqrt{\frac{\sigma_X^2}{\sigma_X^2 + \sigma_Z^2}}
 \end{aligned}$$

- Thus $h(X, Y) = -E_{p(x,y)}[\log_2 p(x, y)] = \frac{1}{2} \log((2\pi e)^2 \det|\mathbf{K}_{xy}|)$

$$\begin{aligned}
 \det|\mathbf{K}_{xy}| &= \det \begin{vmatrix} \sigma_X^2 & \rho_{XY} \sigma_X \sigma_Y \\ \rho_{XY} \sigma_X \sigma_Y & \sigma_Y^2 \end{vmatrix} \\
 &= \sigma_X^2 \sigma_Y^2 - \rho_{XY}^2 \sigma_X^2 \sigma_Y^2 = (1 - \rho_{XY}^2) \sigma_X^2 \sigma_Y^2 = \sigma_X^2 \sigma_Z^2
 \end{aligned}$$

Continuous R.Vs.: MI for univariate Gaussians

- Substituting

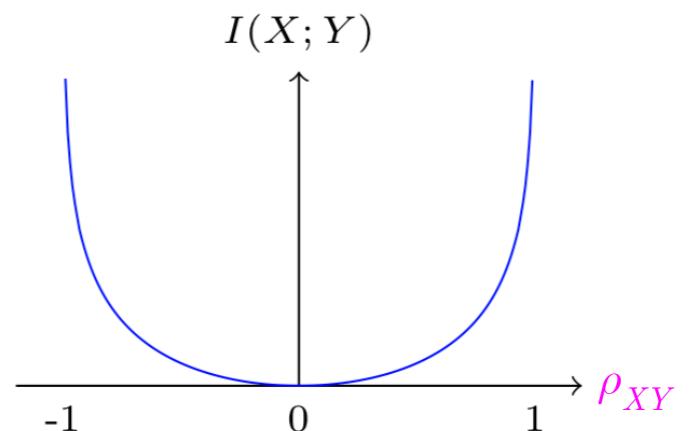
$$\begin{aligned}
 I(X;Y) &= h(X) + h(Y) - h(X,Y) \\
 &= \frac{1}{2} \log(2\pi e \sigma_X^2) + \frac{1}{2} \log \left(2\pi e \underbrace{\sigma_Y^2}_{\sigma_X^2 + \sigma_Z^2} \right) - \frac{1}{2} \log((2\pi e)^2 \sigma_X^2 \sigma_Z^2)
 \end{aligned}$$

$$= \frac{1}{2} \log \left(\frac{(2\pi e)^2}{(2\pi e)^2} \frac{\sigma_X^2 (\sigma_X^2 + \sigma_Z^2)}{\sigma_X^2 \sigma_Z^2} \right) = \frac{1}{2} \log \left(1 + \frac{\sigma_X^2}{\sigma_Z^2} \right) = \frac{1}{2} \log \left(\frac{1}{1 - \rho_{XY}^2} \right)$$

$$\Rightarrow \rho_{XY}^2 = 1 - e^{-2I(X;Y)}$$

$$I(X;Y) = 0 : \rho_{XY}^2 = 1 - e^0 = 0$$

$$I(X;Y) \rightarrow \infty : \rho_{XY}^2 \rightarrow 1$$



Continuous R.Vs.: MI for multivariate Gaussians

- Consider a transformation of Gaussian random variable $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_{\mathbf{xx}})$ by adding an independent random variable $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_{\mathbf{zz}})$, i.e., $\mathbf{Y} = \mathbf{X} + \mathbf{Z}$

$$\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_{\mathbf{xx}}) \longrightarrow \boxed{p(\mathbf{y}|\mathbf{x})} \longrightarrow \mathbf{Y} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_{\mathbf{xx}} + \mathbf{K}_{\mathbf{zz}})$$

$$\mathbf{Y} = \mathbf{X} + \mathbf{Z}$$

$$\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_{\mathbf{zz}})$$

- Compute mutual information via entropy and show the symmetry

$$I(\mathbf{X}; \mathbf{Y}) = h(\mathbf{Y}) - h(\mathbf{Y}|\mathbf{X}) = \frac{1}{2} \log \left(\frac{\det |\mathbf{K}_{\mathbf{xx}} + \mathbf{K}_{\mathbf{zz}}|}{\det |\mathbf{K}_{\mathbf{zz}}|} \right)$$

$$I(\mathbf{X}; \mathbf{Y}) = h(\mathbf{X}) - h(\mathbf{X}|\mathbf{Y}) = \frac{1}{2} \log \left(\frac{\det |\mathbf{K}_{\mathbf{xx}} + \mathbf{K}_{\mathbf{zz}}|}{\det |\mathbf{K}_{\mathbf{zz}}|} \right)$$

Continuous R.Vs.: MI for multivariate Gaussians

$$I(\mathbf{X}; \mathbf{Y}) = h(\mathbf{Y}) - h(\mathbf{Y} | \mathbf{X}) = \frac{1}{2} \log \left(\frac{\det |\mathbf{K}_{\mathbf{xx}} + \mathbf{K}_{\mathbf{zz}}|}{\det |\mathbf{K}_{\mathbf{zz}}|} \right)$$

- Proof

- Random variable $\mathbf{Y} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}_{\mathbf{xx}} + \mathbf{K}_{\mathbf{zz}})$

$$h(\mathbf{Y}) = \frac{1}{2} \ln \left((2\pi e)^N \det |\mathbf{K}_{\mathbf{xx}} + \mathbf{K}_{\mathbf{zz}}| \right)$$

- Since $\mathbf{Y} = \mathbf{X} + \mathbf{Z}$

$$h(\mathbf{Y} | \mathbf{X}) = h(\mathbf{X} + \mathbf{Z} | \mathbf{X}) = h(\mathbf{Z}) = \frac{1}{2} \log \left((2\pi e)^N \det |\mathbf{K}_{\mathbf{zz}}| \right)$$

- Finally

$$I(\mathbf{X}; \mathbf{Y}) = h(\mathbf{Y}) - h(\mathbf{Y} | \mathbf{X}) = \frac{1}{2} \log \left(\frac{(2\pi e)^N \det |\mathbf{K}_{\mathbf{xx}} + \mathbf{K}_{\mathbf{zz}}|}{(2\pi e)^N \det |\mathbf{K}_{\mathbf{zz}}|} \right)$$

Continuous R.Vs.: MI for multivariate Gaussians

- Practical simplification

$$I(\mathbf{X}; \mathbf{Y}) = h(\mathbf{Y}) - h(\mathbf{Y} | \mathbf{X}) = \frac{1}{2} \log \left(\frac{(2\pi e)^N \det |\mathbf{K}_{\mathbf{xx}} + \mathbf{K}_{\mathbf{zz}}|}{(2\pi e)^N \det |\mathbf{K}_{\mathbf{zz}}|} \right)$$

- The term $\det |\mathbf{K}_{\mathbf{zz}}|$

$$\mathbf{K}_{\mathbf{zz}} = \mathbf{U}\Sigma_{\mathbf{zz}}\mathbf{U}^T, \text{ with } \Sigma_{\mathbf{zz}} = \text{diag}(\sigma_1, \dots, \sigma_N) \text{ or } \Sigma_{ii} = \sigma_i, \forall i \text{ and } \mathbf{U}\mathbf{U}^T = \mathbf{I}$$

$$\det |\mathbf{K}_{\mathbf{zz}}| = \left| \det \mathbf{U}\Sigma_{\mathbf{zz}}\mathbf{U}^T \right| = \underbrace{\left| \det \mathbf{U}\mathbf{U}^T \right|}_{\mathbf{I}} \left| \det \Sigma_{\mathbf{zz}} \right| = \left| \det \Sigma_{\mathbf{zz}} \right| = \prod_{i=1}^N \sigma_i$$

- The term $\det |\mathbf{K}_{\mathbf{xx}} + \mathbf{K}_{\mathbf{zz}}|$

$$\begin{aligned} \det |\mathbf{K}_{\mathbf{xx}} + \mathbf{K}_{\mathbf{zz}}| &= \det |\mathbf{K}_{\mathbf{xx}} + \mathbf{U}\Sigma_{\mathbf{zz}}\mathbf{U}^T| = \det |\mathbf{U}| v \det |\mathbf{U}^T| \\ &= \det \left| \underbrace{\mathbf{U}^T \mathbf{K}_{\mathbf{xx}} \mathbf{U}}_{\mathbf{A}} + \Sigma_{\mathbf{zz}} \right| = \det |\mathbf{A} + \Sigma_{\mathbf{zz}}| \leq \prod_{i=1}^N (\mathbf{A}_{ii} + \sigma_i) \end{aligned}$$

Continuous R.Vs.: MI for multivariate Gaussians

- Practical simplification

$$\begin{aligned} I(\mathbf{X}; \mathbf{Y}) &= h(\mathbf{Y}) - h(\mathbf{Y} | \mathbf{X}) = \frac{1}{2} \log \left\{ \frac{(2\pi e)^N \det |\mathbf{K}_{\mathbf{xx}} + \mathbf{K}_{\mathbf{zz}}|}{(2\pi e)^N \det |\mathbf{K}_{\mathbf{zz}}|} \right\} \\ &\leq \frac{1}{2} \log \left\{ \frac{\prod_{i=1}^N (\mathbf{A}_{ii} + \sigma_i)}{\prod_{i=1}^N \sigma_i} \right\}, \text{ where } \mathbf{A} = \mathbf{U}^T \mathbf{K}_{\mathbf{xx}} \mathbf{U} \\ &= \frac{1}{2} \log \left\{ \prod_{i=1}^N \frac{(\mathbf{A}_{ii} + \sigma_i)}{\sigma_i} \right\} \\ &= \frac{1}{2} \log \left\{ \prod_{i=1}^N \left(1 + \frac{\mathbf{A}_{ii}}{\sigma_i} \right) \right\} \\ &= \frac{1}{2} \sum_{i=1}^N \log \left(1 + \frac{\mathbf{A}_{ii}}{\sigma_i} \right) \end{aligned}$$

Scope

- Information theoretic measures for discrete/continuous variables
 - Entropy
 - Conditional entropy
 - Joint entropy
 - Relative entropy (KL-divergence)
 - Cross entropy
 - Mutual information
 - Additional topics:
 - f-divergence
 - practical computations from samples
 - Estimation of KL- and f-divergences by density ratio trick
 - Variational approximation of mutual information

f-divergences

□
Definition (f-divergence): f-divergences or Ali-Silvey distances between $p(x)$ and $q(x)$

$$D_f(p \parallel q) = \int_{\mathcal{X}} p(x)f\left(\frac{q(x)}{p(x)}\right)dx$$

where the function (a.k.a. generator function) $f : \mathbb{R}_+ \rightarrow \mathbb{R}$ is convex and satisfy
 $f(1) = 0$

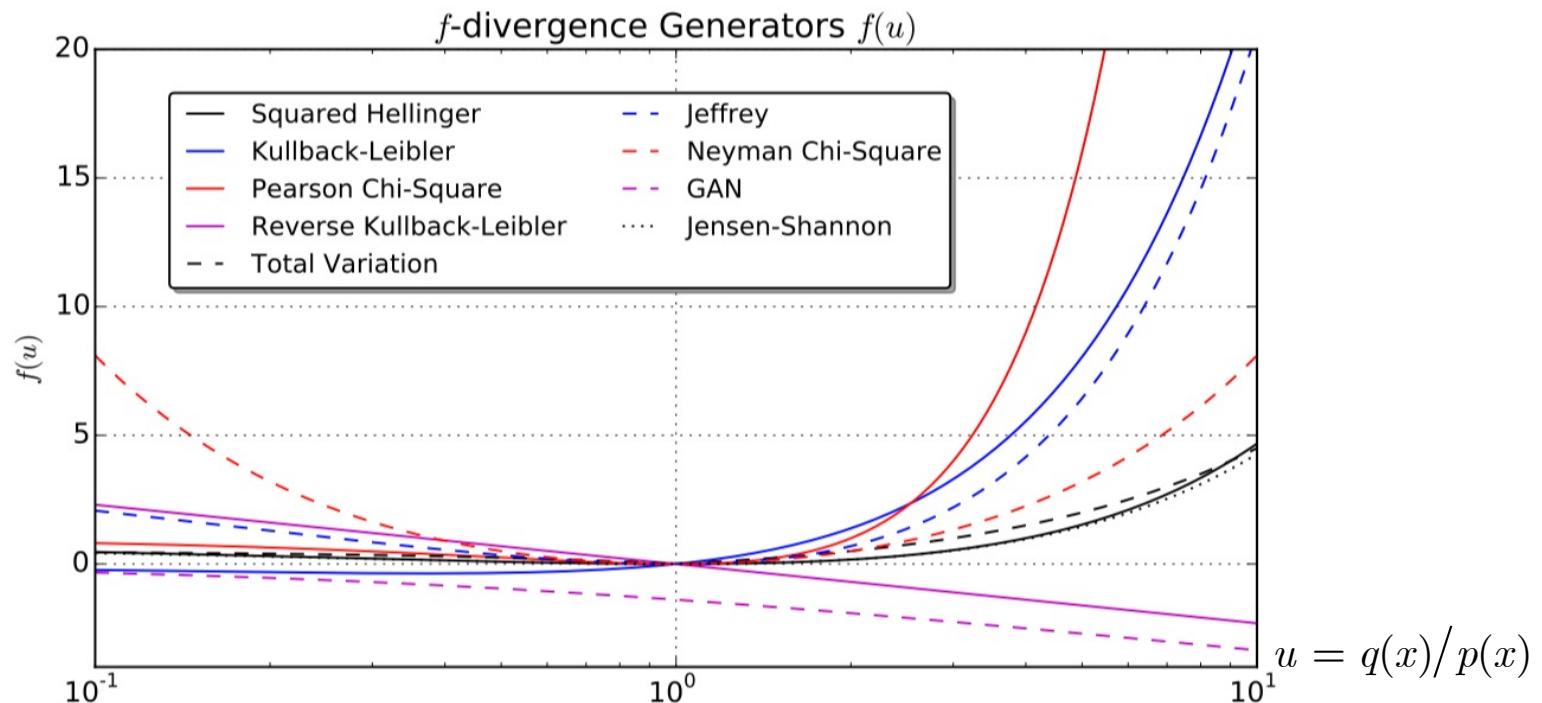
- Recall: KLD

$$D_{\text{KL}}(p \parallel q) = \int_{-\infty}^{\infty} p(x) \log\left(\frac{p(x)}{q(x)}\right) dx = E_{p(x)} \left[\log_2 \frac{p(x)}{q(x)} \right]$$

- Remark: many works define $D_f(p \parallel q) = \int_{\mathcal{X}} q(x)f\left(\frac{p(x)}{q(x)}\right)dx$

f-divergences

- Different choices of $f(u)$



<https://arxiv.org/pdf/1606.00709.pdf>

f-divergences

- Different choices of $f(u)$

Name of the f -divergence	Formula $D_f(p \parallel q)$	Generator $f(u)$ with $f(1) = 0$
Total variation (metric)	$\frac{1}{2} \int p(x) - q(x) d\nu(x)$	$\frac{1}{2} u - 1 $
Squared Hellinger	$\int (\sqrt{p(x)} - \sqrt{q(x)})^2 d\nu(x)$	$(\sqrt{u} - 1)^2$
Pearson χ_P^2	$\int \frac{(q(x) - p(x))^2}{p(x)} d\nu(x)$	$(u - 1)^2$
Neyman χ_N^2	$\int \frac{(p(x) - q(x))^2}{q(x)} d\nu(x)$	$\frac{(1-u)^2}{u}$
Pearson-Vajda χ_P^k	$\int \frac{(q(x) - \lambda p(x))^k}{p^{k-1}(x)} d\nu(x)$	$(u - 1)^k$
Pearson-Vajda $ \chi _P^k$	$\int \frac{ q(x) - \lambda p(x) ^k}{p^{k-1}(x)} d\nu(x)$	$ u - 1 ^k$
Kullback-Leibler	$\int p(x) \log \frac{p(x)}{q(x)} d\nu(x)$	$-\log u$
reverse Kullback-Leibler	$\int q(x) \log \frac{q(x)}{p(x)} d\nu(x)$	$u \log u$
α -divergence	$\frac{4}{1-\alpha^2} (1 - \int p^{\frac{1-\alpha}{2}}(x) q^{1+\alpha}(x) d\nu(x))$	$\frac{4}{1-\alpha^2} (1 - u^{\frac{1+\alpha}{2}})$
Jensen-Shannon	$\frac{1}{2} \int (p(x) \log \frac{2p(x)}{p(x)+q(x)} + q(x) \log \frac{2q(x)}{p(x)+q(x)}) d\nu(x)$	$-(u+1) \log \frac{1+u}{2} + u \log u$

<https://arxiv.org/pdf/1309.3029.pdf>

f-divergences: TV, Pearson and KLD

- **Examples** $D_f(p \parallel q) = \int_{\mathcal{X}} p(x)f\left(\frac{q(x)}{p(x)}\right)dx$
 - Total variation distance $f(u) = \frac{1}{2}|u - 1|$
$$D_{\text{TV}}(p \parallel q) = \frac{1}{2} \int_{\mathcal{X}} p(x) \left| \frac{q(x)}{p(x)} - 1 \right| dx = \frac{1}{2} \int_{\mathcal{X}} |q(x) - p(x)| dx$$
 - Pearson χ^2 – divergence $f(u) = (u - 1)^2$
$$D_{\chi^2}(p \parallel q) = \int_{\mathcal{X}} p(x) \left(\frac{q(x)}{p(x)} - 1 \right)^2 dx = \int_{\mathcal{X}} \frac{(q(x) - p(x))^2}{p(x)} dx$$
 - Kullback-Leibler divergence $f(u) = -\log u$
$$D_{\text{KL}}(p \parallel q) = - \int_{\mathcal{X}} p(x) \log \frac{q(x)}{p(x)} dx = \int_{\mathcal{X}} p(x) \log \frac{p(x)}{q(x)} dx$$

f-divergences: useful inequality

- **TV and KLD**

$$D_{\text{TV}}(p \parallel q)^2 = \|p(x) - q(x)\|_1^2 \leq \frac{1}{2 \log e} D_{\text{KL}}(p \parallel q)$$

- Additionally

$$D_{\text{TV}}(p \parallel q) = \frac{1}{2} \int_{\mathcal{X}} p(x) \left| \frac{q(x)}{p(x)} - 1 \right| dx = \frac{1}{2} \int_{\mathcal{X}} |q(x) - p(x)| dx$$

and

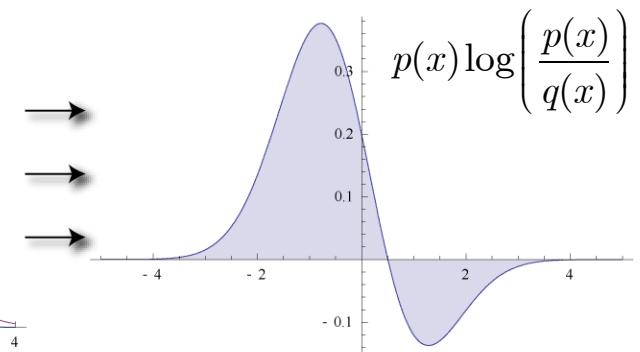
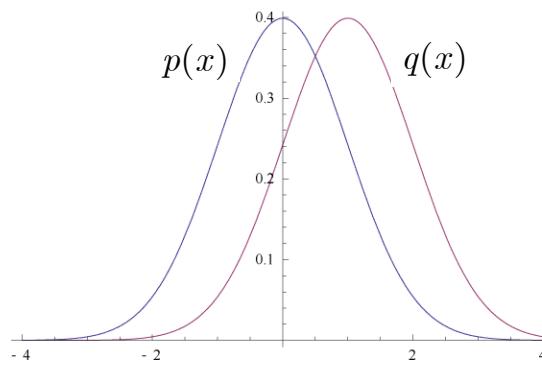
$$\frac{1}{2} |a - b| = \frac{a + b}{2} - \min(a, b) = \max(a, b) - \frac{a + b}{2}$$

$$\begin{aligned} D_{\text{TV}}(p \parallel q) &= \int_{\mathcal{X}} \left(\frac{p(x) + q(x)}{2} - \min(p(x), q(x)) \right) dx \\ &= 1 - \int_{\mathcal{X}} \min(p(x), q(x)) dx = \int_{\mathcal{X}} \max(p(x), q(x)) dx - 1 \end{aligned}$$

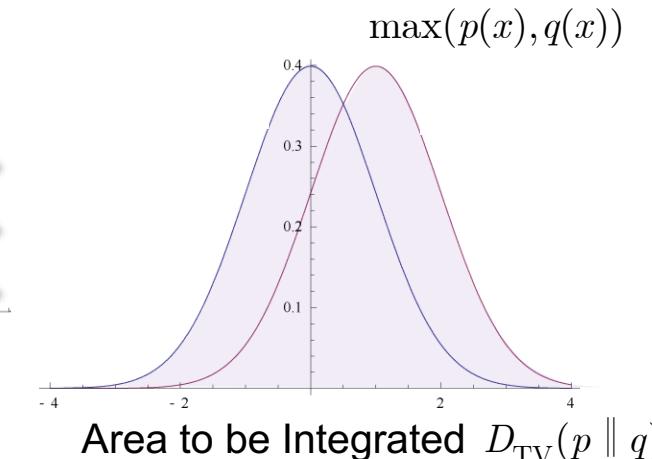
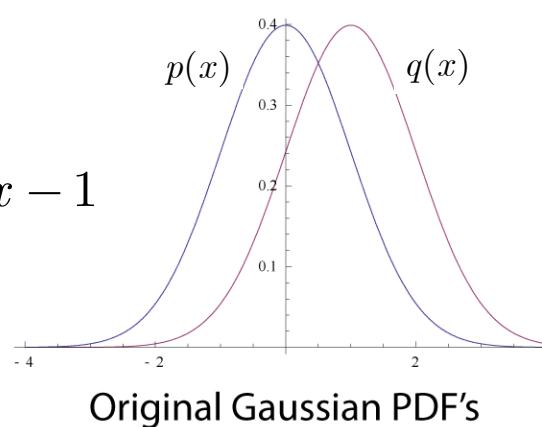
f-divergences: useful inequality

- TV and KLD

$$D_{\text{KL}}(p \parallel q) = \int_{\mathcal{X}} p(x) \log \frac{p(x)}{q(x)} dx$$



$$D_{\text{TV}}(p \parallel q) = \int_{\mathcal{X}} \max(p(x), q(x)) dx - 1$$



f-divergences: useful inequality

- Also a useful inequality

$$\begin{aligned} D_{\text{TV}}(p \parallel q) &= \int_{\mathcal{X}} \left(\frac{p(x) + q(x)}{2} - \min(p(x), q(x)) \right) dx \\ &= 1 - \int_{\mathcal{X}} \min(p(x), q(x)) dx = \int_{\mathcal{X}} \max(p(x), q(x)) dx - 1 \end{aligned}$$

$$\min(a, b) \leq a^\alpha b^{1-\alpha}, \forall a, b > 0$$

$$\min(p(x), q(x)) \leq p(x)^\alpha q(x)^{1-\alpha}$$

$$D_{\text{TV}}(p \parallel q) = 1 - \int_{\mathcal{X}} \min(p(x), q(x)) dx \geq 1 - \int_{\mathcal{X}} p(x)^\alpha q(x)^{1-\alpha} dx$$

f-divergences: Jensen-Shannon

- **Examples**

- Jensen-Shannon divergence $f(u) = -(u + 1)\log \frac{1+u}{2} + u \log u$

$$D_{\text{JS}}(p \parallel q) = \frac{1}{2} \int \left(p(x) \log \frac{2p(x)}{p(x) + q(x)} + q(x) \log \frac{2q(x)}{p(x) + q(x)} \right) dx$$

- Alternative definition: **symmetric KLD**

$$D_{\text{JS}}(p \parallel q) = \frac{1}{2} D_{\text{KL}}(p \parallel m) + \frac{1}{2} D_{\text{KL}}(q \parallel m) \quad m(x) = \frac{1}{2}(p(x) + q(x))$$

$$D_{\text{KL}}(p \parallel m) = \int_{\mathcal{X}} p(x) \log \frac{2p(x)}{p(x) + q(x)} dx$$

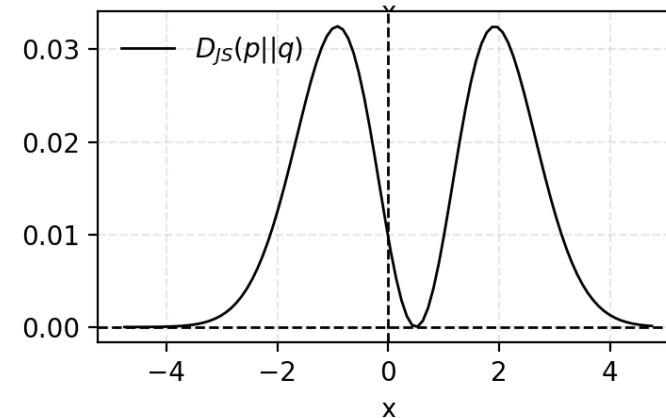
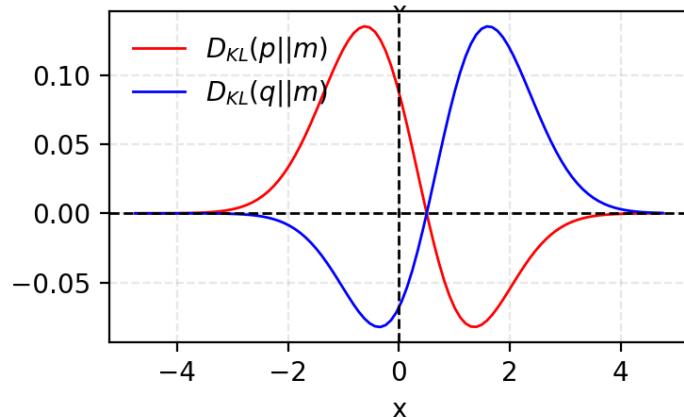
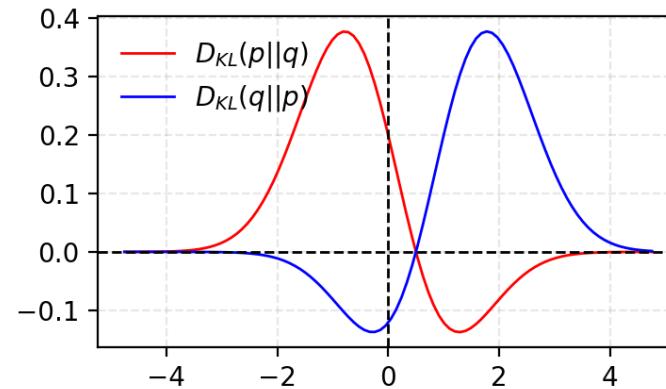
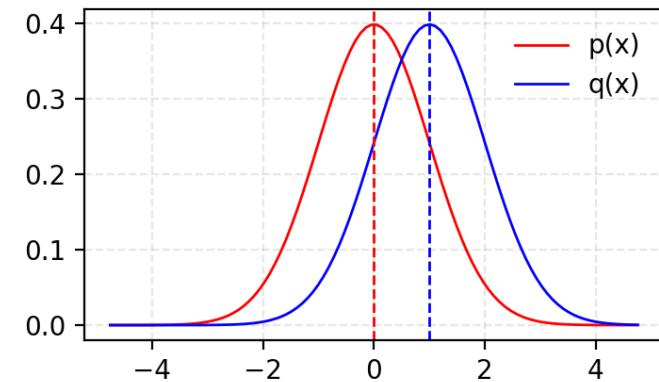
$$D_{\text{KL}}(q \parallel m) = \int_{\mathcal{X}} q(x) \log \frac{2q(x)}{p(x) + q(x)} dx$$

f-divergences: Jensen-Shannon

$$p(x) = \mathcal{N}(0,1)$$

$$q(x) = \mathcal{N}(1,1)$$

$$m(x) = \frac{1}{2} (p(x) + q(x))$$



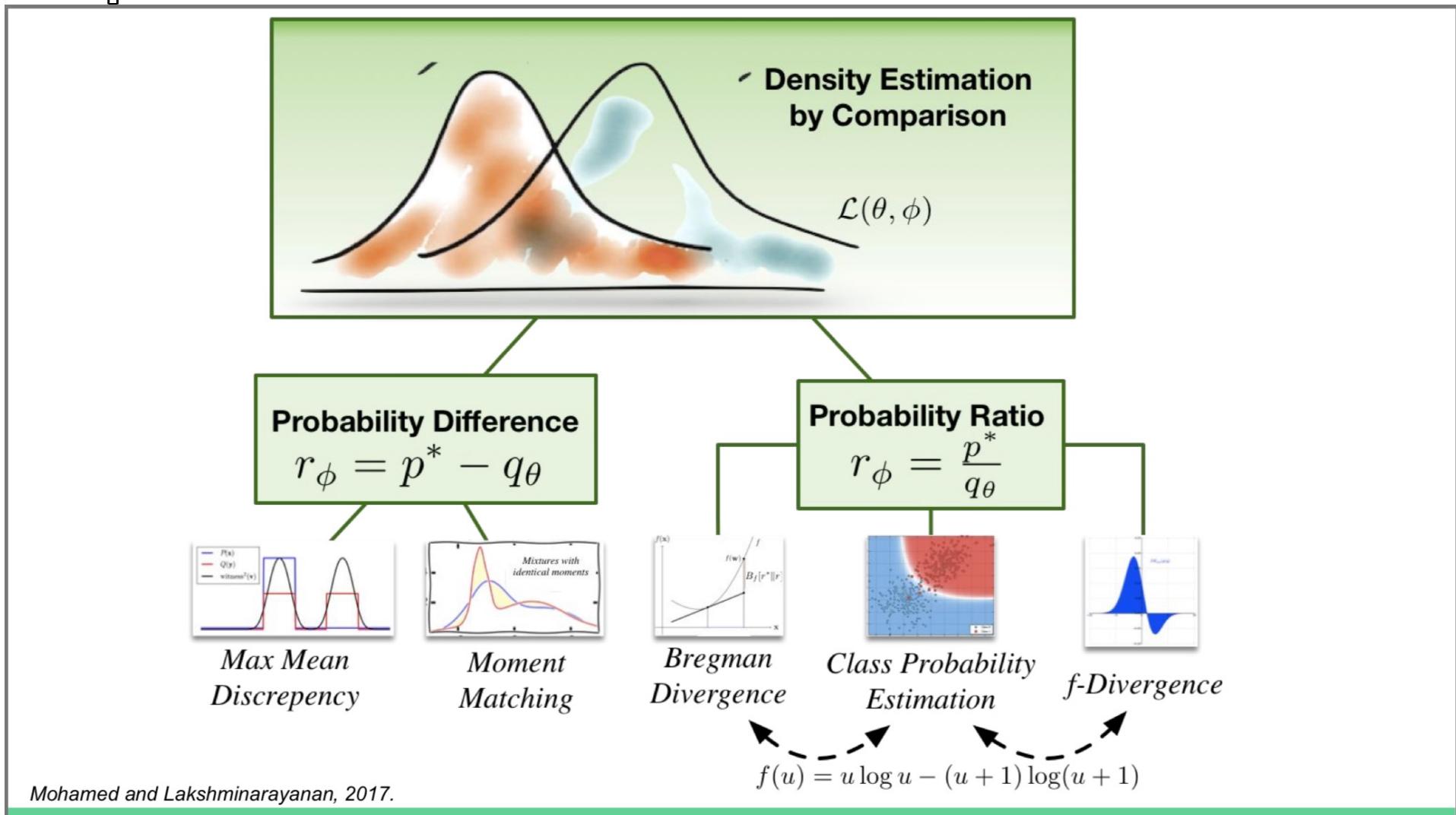
- KL divergence is **asymmetric** while JS divergence is **symmetric**

<https://lilianweng.github.io/lil-log/2017/08/20/from-GAN-to-WGAN.html>

Scope

- Information theoretic measures for discrete/continuous variables
 - Entropy
 - Conditional entropy
 - Joint entropy
 - Relative entropy (KL-divergence)
 - Cross entropy
 - Mutual information
 - Additional topics:
 - f-divergence
 - practical computations from samples
 - Estimation of KL- and f-divergences by density ratio trick
 - Variational approximation of mutual information

f-divergence



<https://www.shakirm.com/slides/DeepGenModelsTutorial.pdf>

Estimation of density ratio for KL divergence

- KL divergence between distributions $p(x)$ and $q(x)$ is defined as

▫

$$D_{\text{KL}}(p \parallel q) = \int_{\mathcal{X}} p(x) \log \left(\frac{p(x)}{q(x)} \right) dx$$

- It can be expressed in terms of density ratio $r(x) = \frac{p(x)}{q(x)}$

$$D_{\text{KL}}(p \parallel q) = \int_{\mathcal{X}} p(x) \log r(x) dx = E_{p(x)} [\log r(x)]$$

- Remark: similar situation for f-divergence

$$D_f(p \parallel q) = \int_{\mathcal{X}} p(x) f \left(\frac{q(x)}{p(x)} \right) dx = \int_{\mathcal{X}} p(x) f \left(r^{-1}(x) \right) dx = E_{p(x)} [f(r^{-1}(x))]$$

Estimation of density ratio for KL divergence

- **Challenges**

- In both KL- or f-divergence cases

$$D_{\text{KL}}(p \parallel q) = E_{p(x)}[\log r(x)] \quad \text{and} \quad D_f(p \parallel q) = E_{p(x)}[f(r^{-1}(x))]$$

we need to know how to compute

- **Expectation** $E_{p(x)}[\cdot]$ wrt $p(x)$
 - **Density ratio** $r(x) = \frac{p(x)}{q(x)}$

Estimation of density ratio for KL divergence

- Practical issues:

Case 1

- Suppose $p(x)$ and $q(x)$ **are known** (explicitly) but
 - It is difficult to compute mathematically the expectation $E_{p(x)}[.]$ but still possible to compute the ratio $r(x) = \frac{p(x)}{q(x)}$
- It is intractable to compute the ratio $r(x) = \frac{p(x)}{q(x)}$

Case 2

- Suppose $p(x)$ and $q(x)$ are not known explicitly but we know some samples generated from these distributions (univariate or multivariate)

$$x_p^{(1)}, \dots, x_p^{(N_p)} \sim p(x) \quad x_q^{(1)}, \dots, x_q^{(N_q)} \sim q(x)$$

implicit distributions

- How to estimate $D_{\text{KL}}(p \parallel q)$ or $D_f(p \parallel q)$ in these cases?

Estimation of density ratio for KL divergence

- Case 1:

- Monte Carlo estimation

$$E_{p(x)}[\cdot] \rightarrow \frac{1}{N} \sum_{i=1}^M$$

- KL-divergence with analytically tractable density ratio $r(x) = \frac{p(x)}{q(x)}$

$$D_{\text{KL}}(p \parallel q) = E_{p(x)}[\log r(x)]$$

$$\approx \frac{1}{M} \sum_{i=1}^M \log r(x_p^{(i)}), x_p^{(i)} \sim p(x)$$

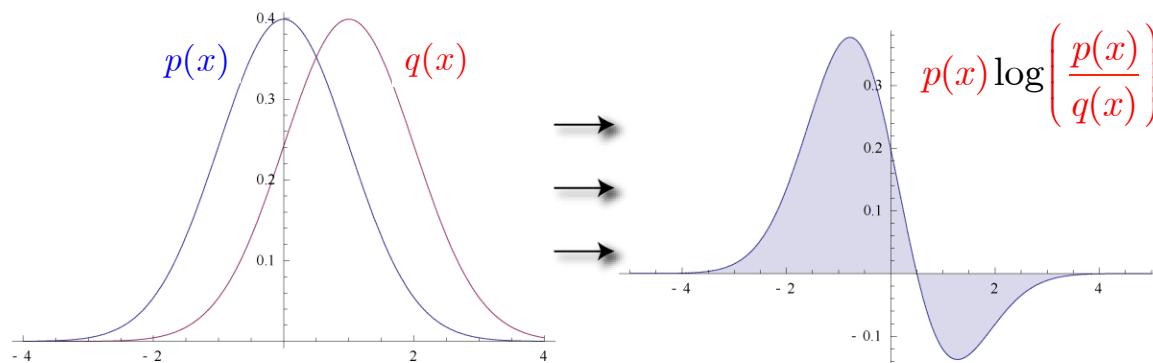
- f-divergence with analytically tractable density ratio $r(x) = \frac{p(x)}{q(x)}$

$$D_f(p \parallel q) = E_{p(x)} \left[\log(r(x))^{-1} \right]$$

$$\approx \frac{1}{M} \sum_{i=1}^M \log(r(x_p^{(i)}))^{-1}, x_p^{(i)} \sim p(x)$$

Estimation of density ratio for KL divergence

- Case 2:
 - Density ratio estimation – implicit distributions
 - It is intractable to compute the ratio $r(x) = \frac{p(x)}{q(x)}$
 - Distributions are unknown
 - Note: we do not need to know distributions explicitly but we only need to know:
 - either their ratio $r(x) = \frac{p(x)}{q(x)}$ like in KL or JS cases
 - or difference $p(x) - q(x)$ like in Pearson or TV cases



Estimation of density ratio for KL divergence

- Case 2:

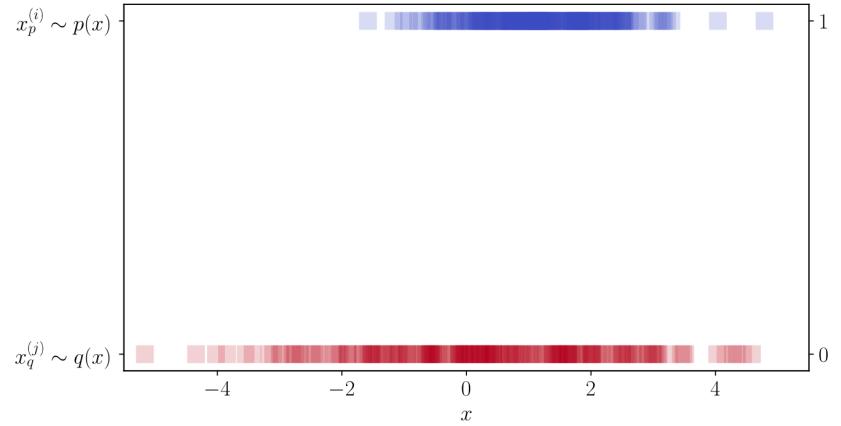
- Suppose we have samples from both distributions

$$x_p^{(1)}, \dots, x_p^{(N_p)} \sim p(x) \quad x_q^{(1)}, \dots, x_q^{(N_q)} \sim q(x)$$

- Then, we form the dataset $\{x_n, y_n\}_{n=1}^N$, where $N = N_p + N_q$
 - We label samples drawn from $p(x)$ as 1 and those drawn from $q(x)$ as 0

$$(x_1, \dots, x_N) = (x_p^{(1)}, \dots, x_p^{(N_p)}, x_q^{(1)}, \dots, x_q^{(N_q)})$$

$$(y_1, \dots, y_N) = (\underbrace{1, \dots, 1}_{N_p}, \underbrace{0, \dots, 0}_{N_q})$$



<https://tiao.io/post/density-ratio-estimation-for-kl-divergence-minimization-between-implicit-distributions/>

Estimation of density ratio for KL divergence

- Case 2:

- By the labeling

$$p(x) = \mathcal{P}(x \mid y = 1) \quad q(x) = \mathcal{P}(x \mid y = 0)$$

- Using Baye's rule, we can write

$$\mathcal{P}(x \mid y) = \frac{\mathcal{P}(y \mid x)\mathcal{P}(x)}{\mathcal{P}(y)}$$

- The density ratio can be re-expressed as:

$$\begin{aligned} r(x) &= \frac{p(x)}{q(x)} = \frac{\mathcal{P}(x \mid y = 1)}{\mathcal{P}(x \mid y = 0)} \\ &= \left(\frac{\mathcal{P}(y = 1 \mid x)\mathcal{P}(x)}{\mathcal{P}(y = 1)} \right) \left(\frac{\mathcal{P}(y = 0 \mid x)\mathcal{P}(x)}{\mathcal{P}(y = 0)} \right)^{-1} = \frac{\mathcal{P}(y = 0)}{\mathcal{P}(y = 1)} \frac{\mathcal{P}(y = 1 \mid x)}{\mathcal{P}(y = 0 \mid x)} \end{aligned}$$

Estimation of density ratio for KL divergence

- Case 2:

- Approximate the ratio of marginal densities by the ratio of sample sizes

$$\frac{\mathcal{P}(y = 0)}{\mathcal{P}(y = 1)} \approx \frac{N_q}{N_p + N_q} \left(\frac{N_p}{N_p + N_q} \right)^{-1} = \frac{N_q}{N_p}$$

- Let us assume from now on that $N_q = N_p$. Thus

$$r(x) = \frac{p(x)}{q(x)} = \frac{\mathcal{P}(y = 0)}{\mathcal{P}(y = 1)} \frac{\mathcal{P}(y = 1 | x)}{\mathcal{P}(y = 0 | x)} = \frac{\mathcal{P}(y = 1 | x)}{\mathcal{P}(y = 0 | x)}$$

- This can be developed as

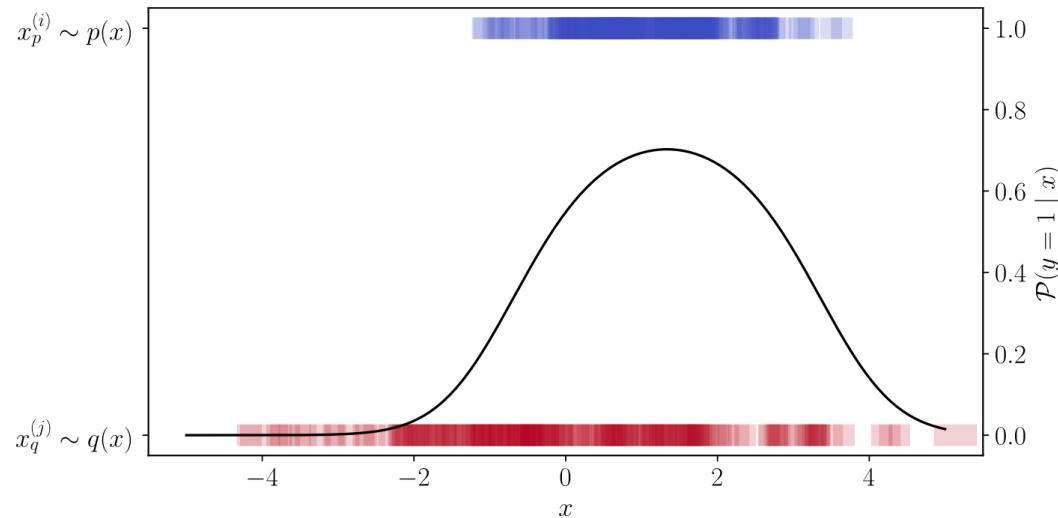
$$\begin{aligned} r(x) &= \frac{\mathcal{P}(y = 1 | x)}{\mathcal{P}(y = 0 | x)} = \frac{\mathcal{P}(y = 1 | x)}{1 - \mathcal{P}(y = 1 | x)} \\ &= \exp \left[\log \frac{\mathcal{P}(y = 1 | x)}{1 - \mathcal{P}(y = 1 | x)} \right] = \exp \left[\sigma^{-1}(\mathcal{P}(y = 1 | x)) \right] \end{aligned}$$

where $\sigma^{-1}(t) = \log(t/(q-t))$ is the *logit* function, or inverse sigmoid function

Estimation of density ratio for KL divergence

- Case 2:
- **Recovering the Class Probability from the Density Ratio**

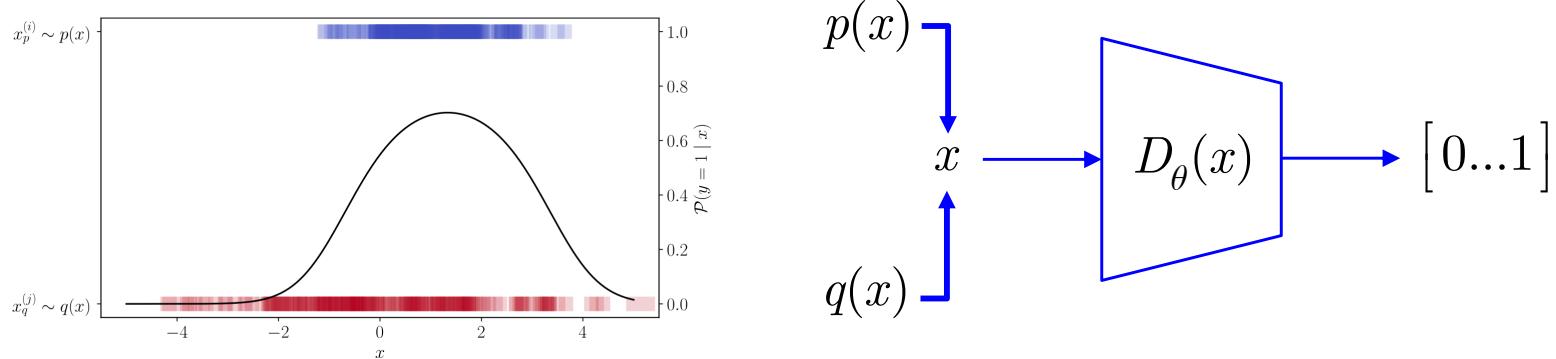
$$\mathcal{P}(y = 1 \mid x) = \sigma(\log r(x)) = \frac{p(x)}{p(x) + q(x)}$$



Note: $\mathcal{P}(y = 1 \mid x)$ is influenced by both $p(x)$ and $q(x)$ that is a source of potential ambiguity

Estimation of density ratio for KL divergence

- Case 2:
- Probabilistic Classification with Logistic Regression
- We can approximate $\mathcal{P}(y = 1 | x)$ by a parametric function $D_\theta(x)$ with parameters θ



- The corresponding density ratio estimator

$$r_\theta(x) = \exp\left[\sigma^{-1}(D_\theta(x))\right] \approx \exp\left[\sigma^{-1}(\mathcal{P}(y = 1 | x))\right] = r(x)$$

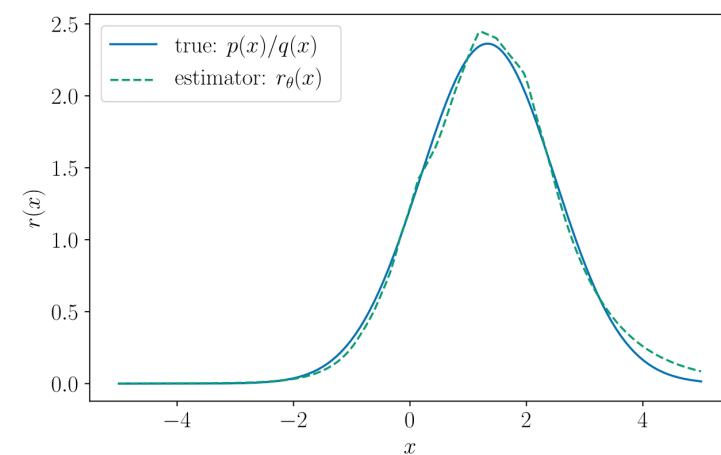
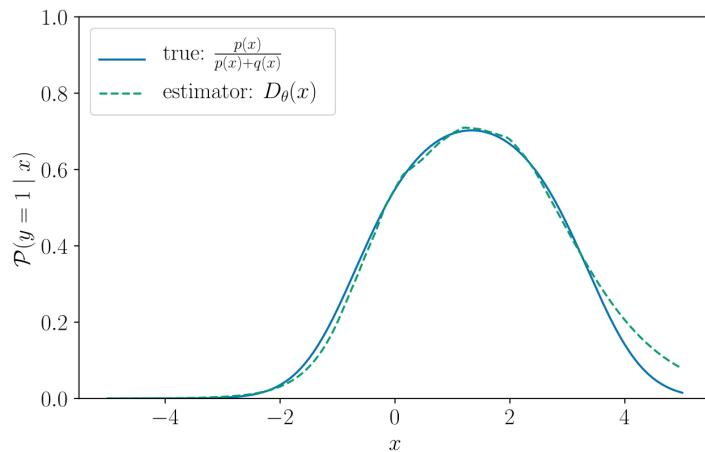
- or alternatively

$$D_\theta(x) = \sigma(\log r_\theta(x)) \approx \sigma(\log r(x)) = \mathcal{P}(y = 1 | x)$$

Estimation of density ratio for KL divergence

- Case 2:
- Training the parameters of estimator $D_\theta(x)$
 - Using binary cross-entropy, one can try to match the score $D_\theta(x)$ with $p(x)$ and the score $(1 - D_\theta(x))$ with $q(x)$

$$\begin{aligned}\mathcal{L}(\theta) &:= -\mathbb{E}_{p(x)}[\log D_\theta(x)] - \mathbb{E}_{q(x)}[\log(1 - D_\theta(x))] \\ &= -\mathbb{E}_{p(x)}[\log \sigma(\log r_\theta(x))] - \mathbb{E}_{q(x)}[\log(1 - \sigma(\log r_\theta(x)))]\end{aligned}$$



Estimation of density ratio for KL divergence

- Case 2:
- Back to Monte Carlo estimation
 - Having obtained an estimate of the log density ratio, it is now feasible to perform Monte Carlo estimation

$$\begin{aligned} D_{\text{KL}}(p \parallel q) &= E_{p(x)}[\log r(x)] \\ &\approx \frac{1}{M} \sum_{i=1}^M \log r(x_p^{(i)}), \quad x_p^{(i)} \sim p(x) \\ &\approx \frac{1}{M} \sum_{i=1}^M \log r_\theta(x_p^{(i)}), \quad x_p^{(i)} \sim p(x) \end{aligned}$$

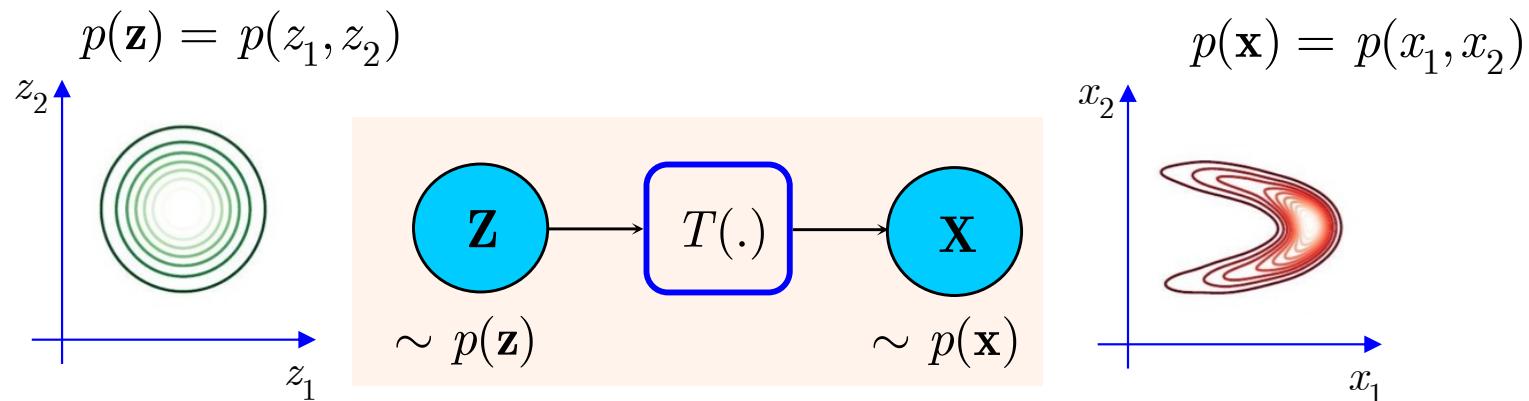
- we draw MC samples from $p(x)$ as before
- but instead of taking the mean of the function $\log r(x_p^{(i)})$ evaluated on these samples (which is unavailable for implicit distributions), we do so on a proxy function $\log r_\theta(x_p^{(i)})$ that is estimated through probabilistic classification as described above.

Scope

- Information theoretic measures for discrete/continuous variables
 - Entropy
 - Conditional entropy
 - Joint entropy
 - Relative entropy (KL-divergence)
 - Cross entropy
 - Mutual information
 - Additional topics:
 - f-divergence
 - practical computations from samples
 - Estimation of KL- and f-divergences by density ratio trick
 - Variational approximation of mutual information

Approximation approach: IT generative models

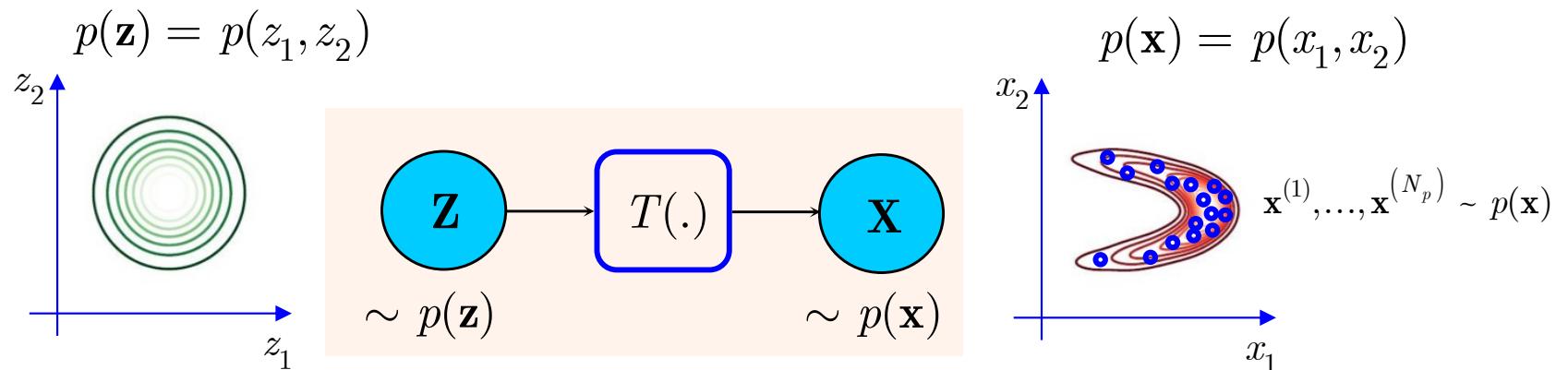
- Suppose we want to build a mapper (generator) that generates samples following the distribution $p(\mathbf{x})$ from the samples generated by the distribution $p(\mathbf{z})$



- If the distributions $p(\mathbf{z})$ and $p(\mathbf{x})$ were known analytically, we would easily compute a mapper (transform) $T(\cdot)$ analytically (see Theme 1)

Approximation approach: IT generative models

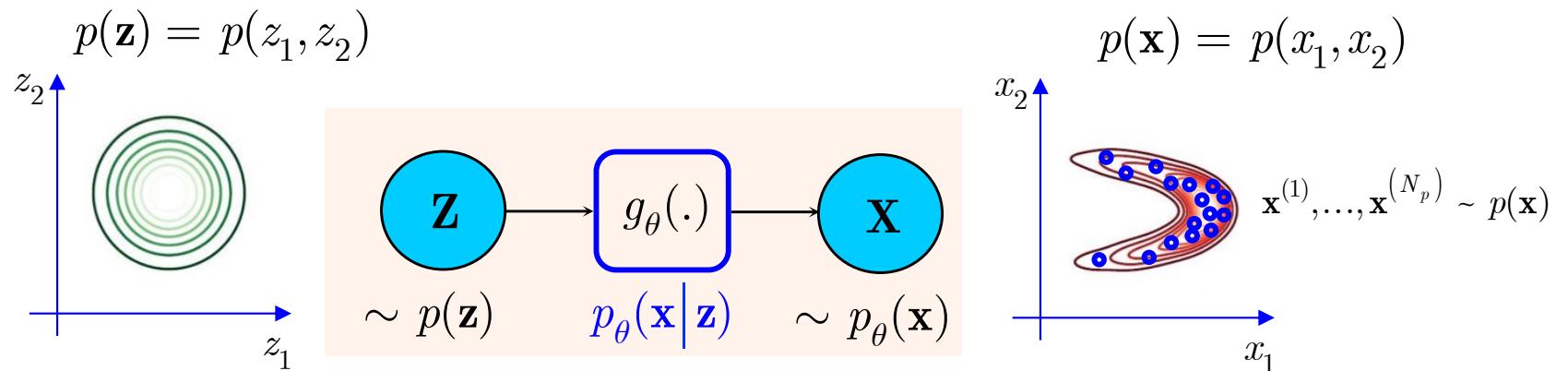
- **In practice:** the distribution $p(\mathbf{z})$ is known and we can easily generate samples from it (no need to parametrize it) while the distribution $p(\mathbf{x})$ is only represented by samples $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N_p)} \sim p(\mathbf{x})$



- Our challenge: to find a transformation $T(\cdot)$ that would map the distribution $p(\mathbf{z})$ into the observed samples $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N_p)} \sim p(\mathbf{x})$
- The “empirical” distribution of observed samples is denoted as $p_D(\mathbf{x})$

Approximation approach: IT generative models

- **Parametrization:** we proceed by parametrizing an unknown transform $T(.)$ by some parametric function $g_\theta(.)$ with a set of parameters θ



- The distribution of generated samples (samples $\mathbf{Z} \sim p(\mathbf{z})$ passed via $p_\theta(\mathbf{x}|\mathbf{z})$) is

$$p_\theta(\mathbf{x}) = \int p_\theta(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

- Our goal would be to find parameters θ , such that $p_\theta(\mathbf{x}) \rightarrow p_D(\mathbf{x})$

Approximation approach: IT generative models

- Possible strategies to find parameters θ such that: $p_\theta(\mathbf{x}) \rightarrow p_D(\mathbf{x})$
 - via KL-divergence minimization

$$\theta = \arg \min_{\theta} D_{\text{KL}}(p_D(\mathbf{x}) \parallel p_\theta(\mathbf{x}))$$

We know how to do it
and also properties of
direct and reverse KLD

- via maximization of mutual information

$$\theta = \arg \max_{\theta} I_\theta(\mathbf{X}; \mathbf{Z})$$

$$\begin{aligned} I_\theta(\mathbf{X}; \mathbf{Z}) &= E_{p(\mathbf{x}, \mathbf{z})} \left[\log \frac{p_\theta(\mathbf{x} | \mathbf{z})}{p_D(\mathbf{x})} \frac{p(\mathbf{z})}{p(\mathbf{z})} \frac{\cancel{p_\theta(\mathbf{x})}}{\cancel{p_\theta(\mathbf{x})}} \right] \\ &= \underbrace{E_{p(\mathbf{x}, \mathbf{z})} \left[\log p_\theta(\mathbf{x} | \mathbf{z}) \right]}_{\propto -\lambda E_{p(\mathbf{x}, \mathbf{z})} \left[\|\mathbf{x} - g_\theta(\mathbf{z})\|_2^2 \right]} - \underbrace{E_{p(\mathbf{x}, \mathbf{z})} \left[\log \frac{p_D(\mathbf{x})}{p_\theta(\mathbf{x})} \right]}_{-D_{\text{KL}}(p_D(\mathbf{x}) \parallel \cancel{p_\theta(\mathbf{x})})} - \underbrace{E_{p(\mathbf{x}, \mathbf{z})} \left[\log \cancel{p_\theta(\mathbf{x})} \right]}_{H(p_D(\mathbf{x}); \cancel{p_\theta(\mathbf{x})})} \\ p(\mathbf{x}, \mathbf{z}) &= p_D(\mathbf{x})p(\mathbf{z}) \end{aligned}$$

Approximation approach: IT generative models

- Possible strategies to find parameters θ such that: $p_\theta(\mathbf{x}) \rightarrow p_D(\mathbf{x})$

$$\theta = \arg \max_{\theta} I_{\theta}(\mathbf{X}; \mathbf{Z})$$

$$= \arg \max_{\theta} \left[-\lambda E_{p(\mathbf{x}, \mathbf{z})} \left[\left\| \mathbf{x} - g_{\theta}(\mathbf{z}) \right\|_2^2 \right] - D_{\text{KL}}(p_D(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})) + H(p_D(\mathbf{x}); p_{\theta}(\mathbf{x})) \right]$$

- Since $H(p_D(\mathbf{x}); p_{\theta}(\mathbf{x})) \geq 0$

$$\theta = \arg \min_{\theta} \left[\lambda E_{p(\mathbf{x}, \mathbf{z})} \left[\left\| \mathbf{x} - g_{\theta}(\mathbf{z}) \right\|_2^2 \right] + D_{\text{KL}}(p_D(\mathbf{x}) \parallel p_{\theta}(\mathbf{x})) \right]$$

We reduced it to two metrics that we know how to compute

- Note a link to the direct minimization of KL- divergence